

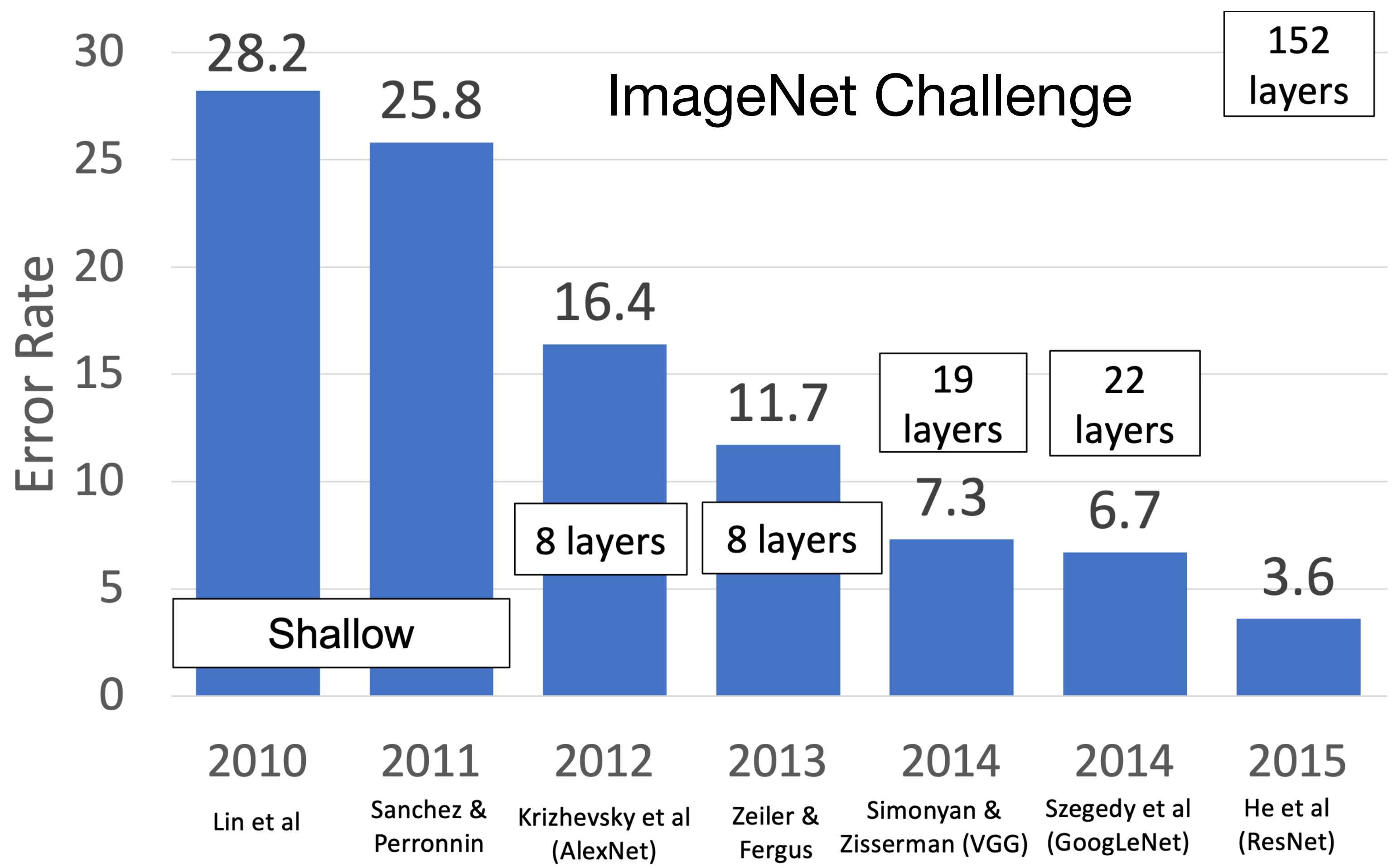


DEEP ROB

Lecture 9
Object Detection
University of Michigan | Department of Robotics



Recap: CNN architectures for Classification



- Some key concepts:
- VGG: stacked Conv+ReLU, pool
- Stem network (GoogLeNet)
- Residual Network (ResNet)



Object Detection: Task definition

Input: Single RGB image

Output: A set of detected objects;

For each object predict:

1. Category label (from a fixed set of labels)
2. Bounding box (four numbers: x, y, width, height)





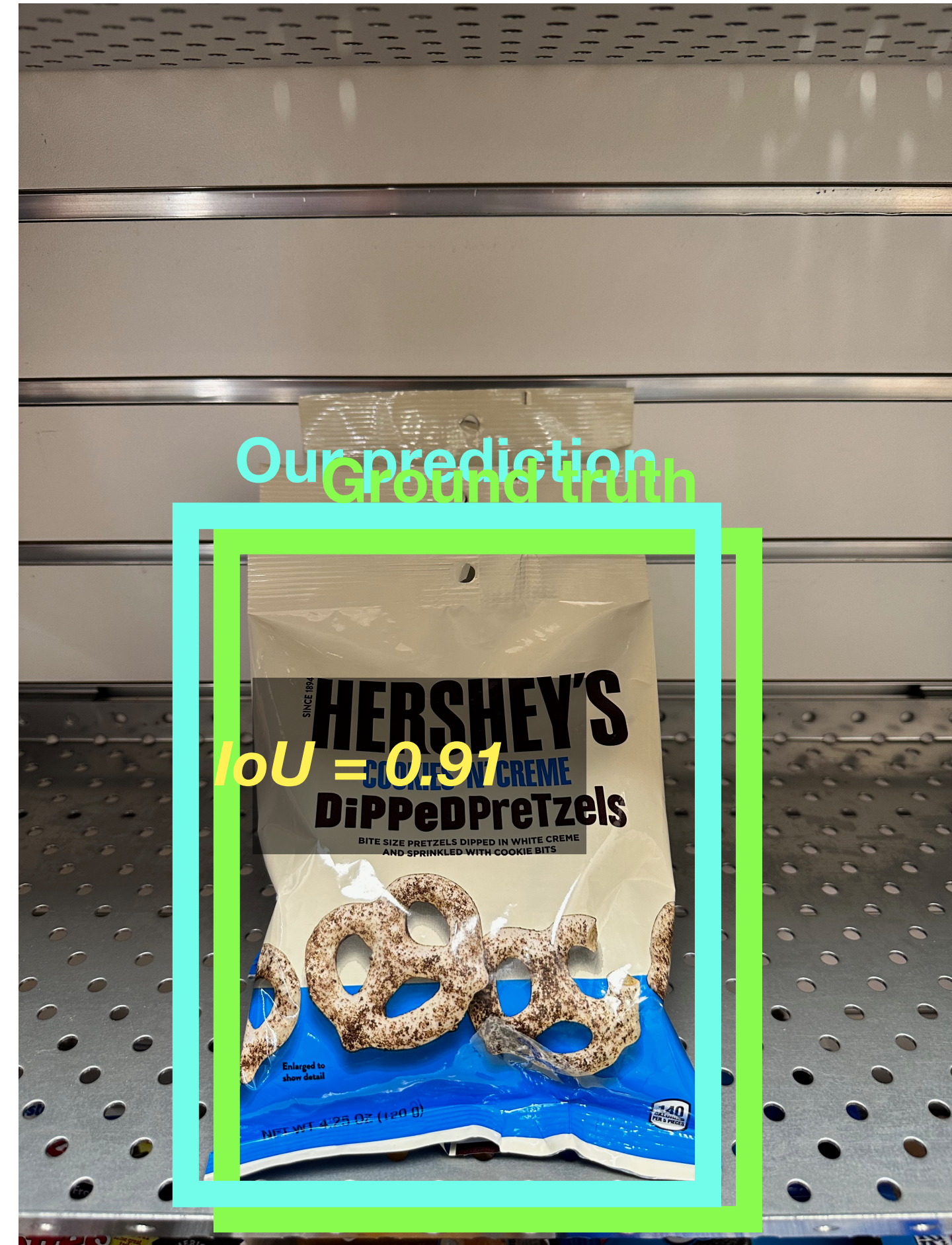
Comparing Boxes: Intersection over Union (IoU)

How can we compare our prediction to the ground-truth box?

Intersection over Union (IoU) (Also called “Jaccard similarity” or “Jaccard index”):

$$\text{IoU} = \frac{\text{Area of Intersection}}{\text{Area of Union}}$$

- IoU > 0.5 is “decent”,
- IoU > 0.7 is “pretty good”,
- IoU > 0.9 is “very good, almost perfect”





Detecting a single object

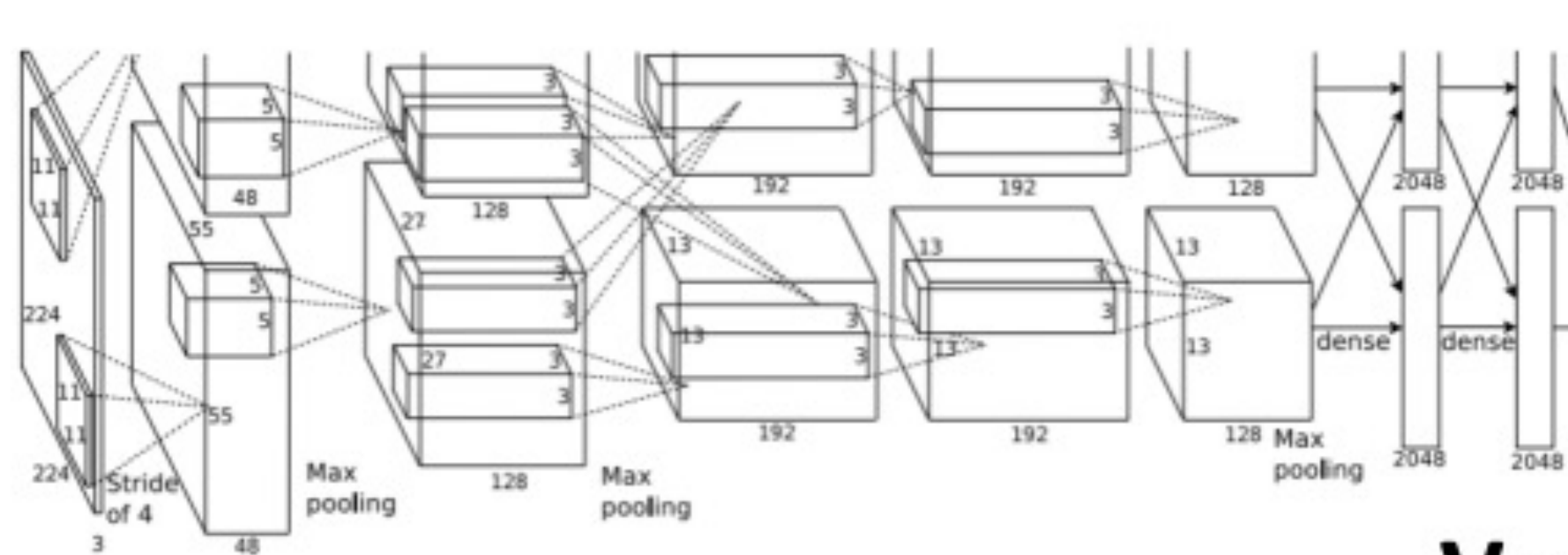


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Fully connected:
4096 to 10

Vector:
4096

What??

Class scores:
 Chocolate Pretzels: 0.9
 Granola Bar: 0.02
 Potato Chips: 0.02
 Water Bottle: 0.02
 Popcorn: 0.01

Correct Label:
Chocolate Pretzels



Treat localization as a regression problem!



Detecting a single object

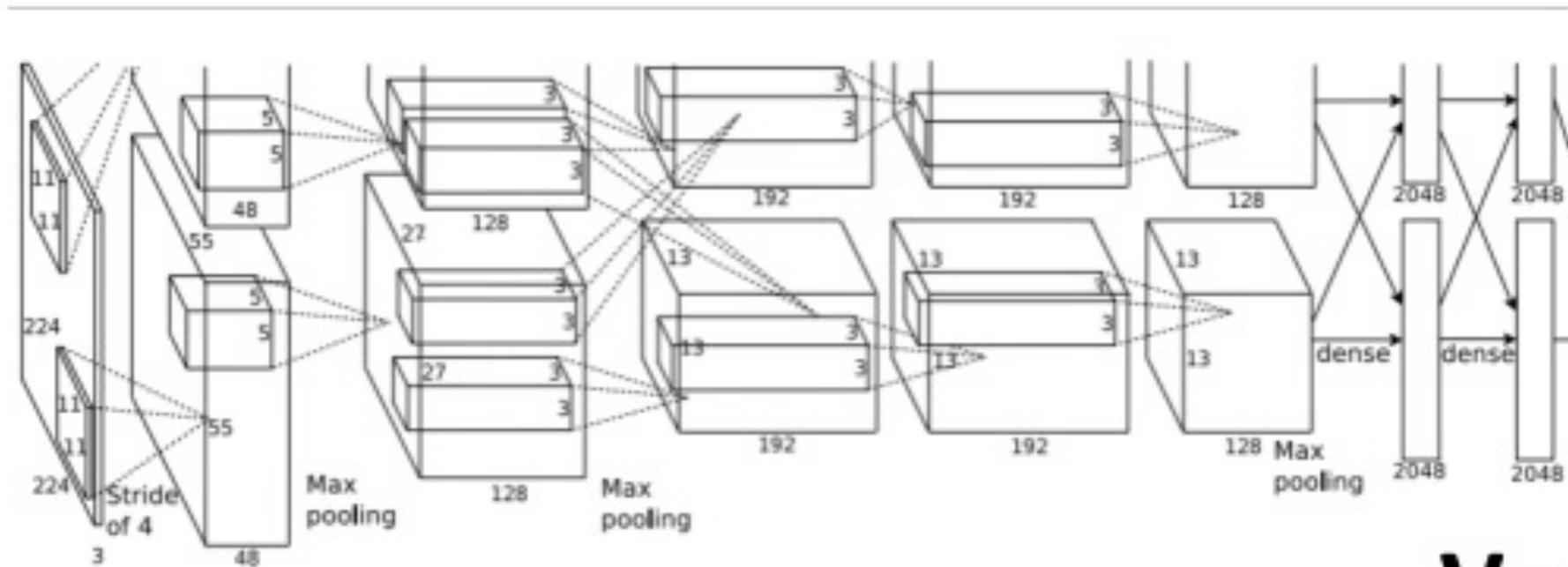


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Vector:
4096

Fully connected:
4096 to 10

What??

Class scores:
 Chocolate Pretzels: 0.9
 Granola Bar: 0.02
 Potato Chips: 0.02
 Water Bottle: 0.02
 Popcorn: 0.01

Correct Label:
Chocolate Pretzels

Softmax Loss

Treat localization as a regression problem!

Fully connected:
4096 to 10

Box coordinates:
(x, y, w, h)

Where??

L2 Loss

Correct coordinates:
(x', y', w', h')



Detecting a single object

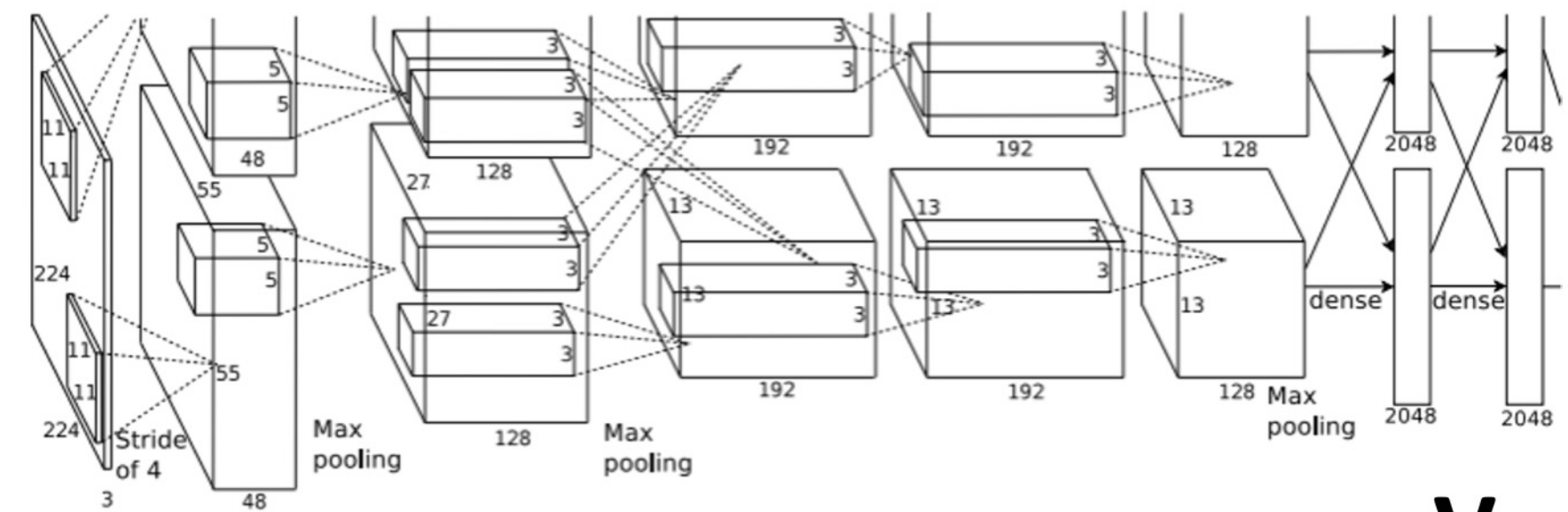


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Fully connected:
4096 to 10

Vector:
4096

Fully connected:
4096 to 10

- Class scores:
- Chocolate Pretzels: 0.9 *What??*
 - Granola Bar: 0.02
 - Potato Chips: 0.02
 - Water Bottle: 0.02
 - Popcorn: 0.01
 -

Correct Label:
Chocolate Pretzels

Softmax Loss

Multitask Loss

Weighted Sum

Loss

$$L = L_{cls} + \lambda L_{reg}$$

L2 Loss

Box coordinates:

(x, y, w, h)

Where??

Correct coordinates:

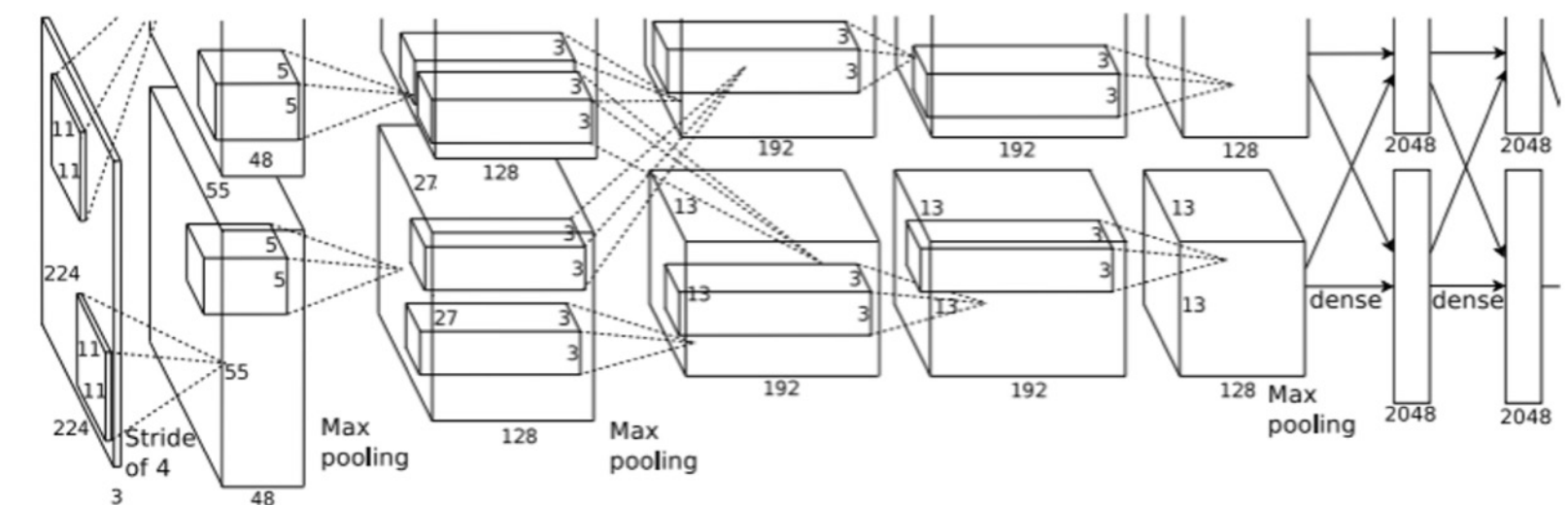
(x', y', w', h')

Treat localization as a regression problem!

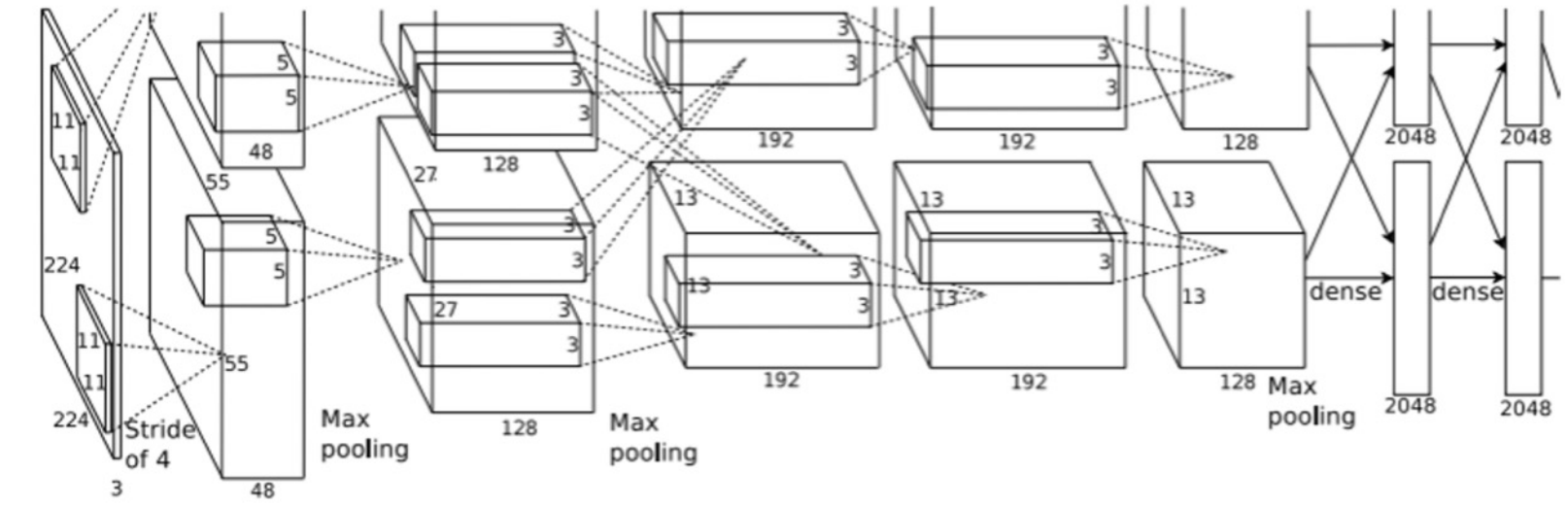




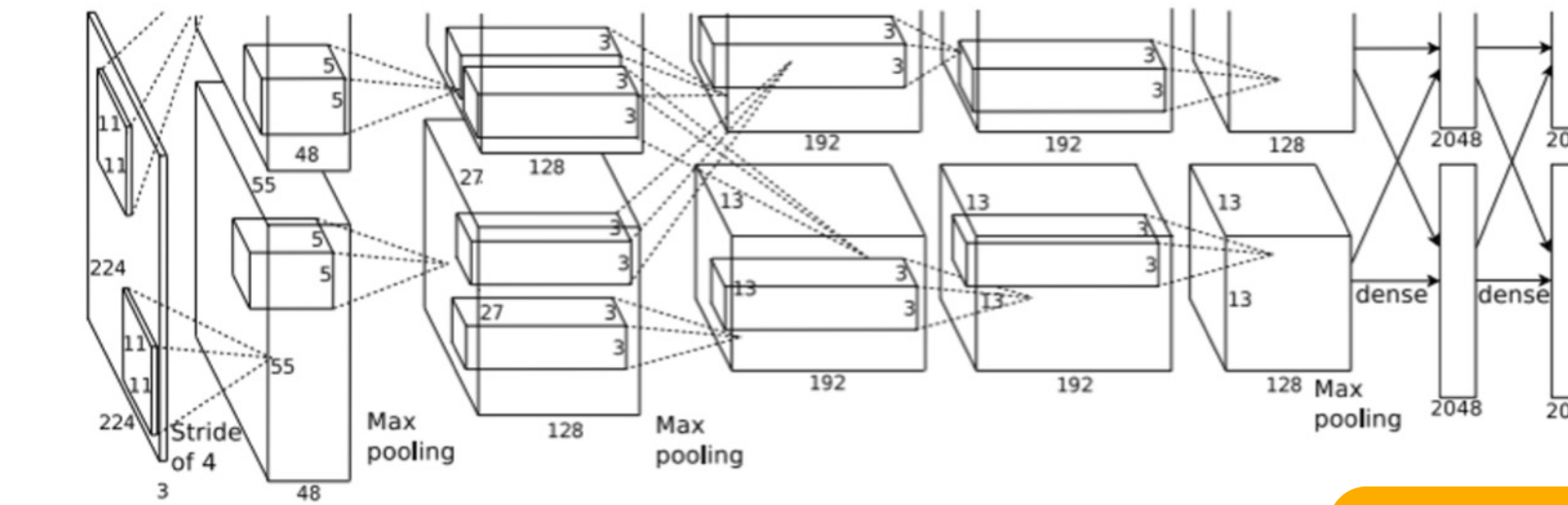
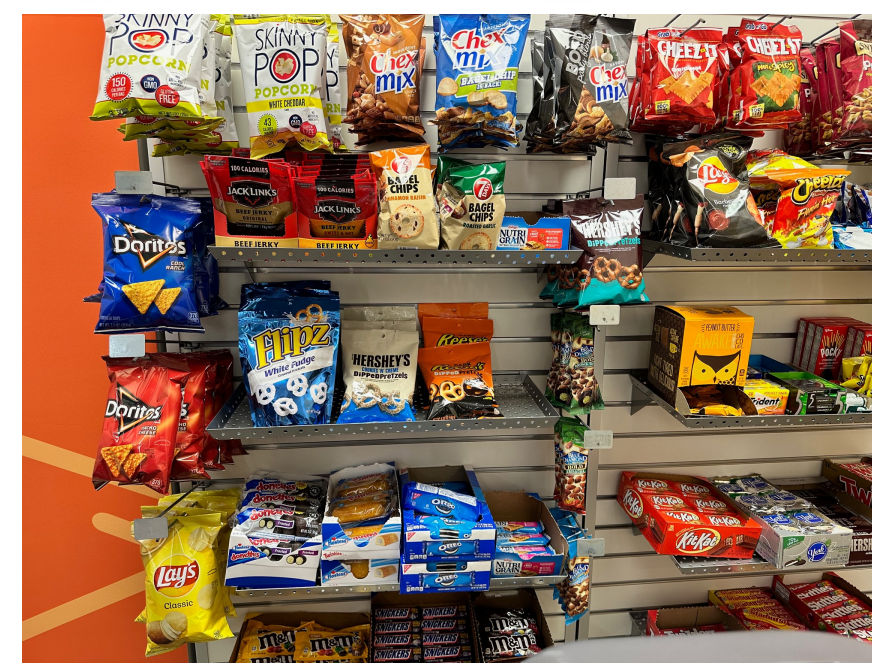
Detecting Multiple Objects



Hershey's: (x, y, w, h)
4 numbers



Hershey's: (x, y, w, h)
Flipz: (x, y, w, h)
Reese's (x, y, w, h)
12 numbers



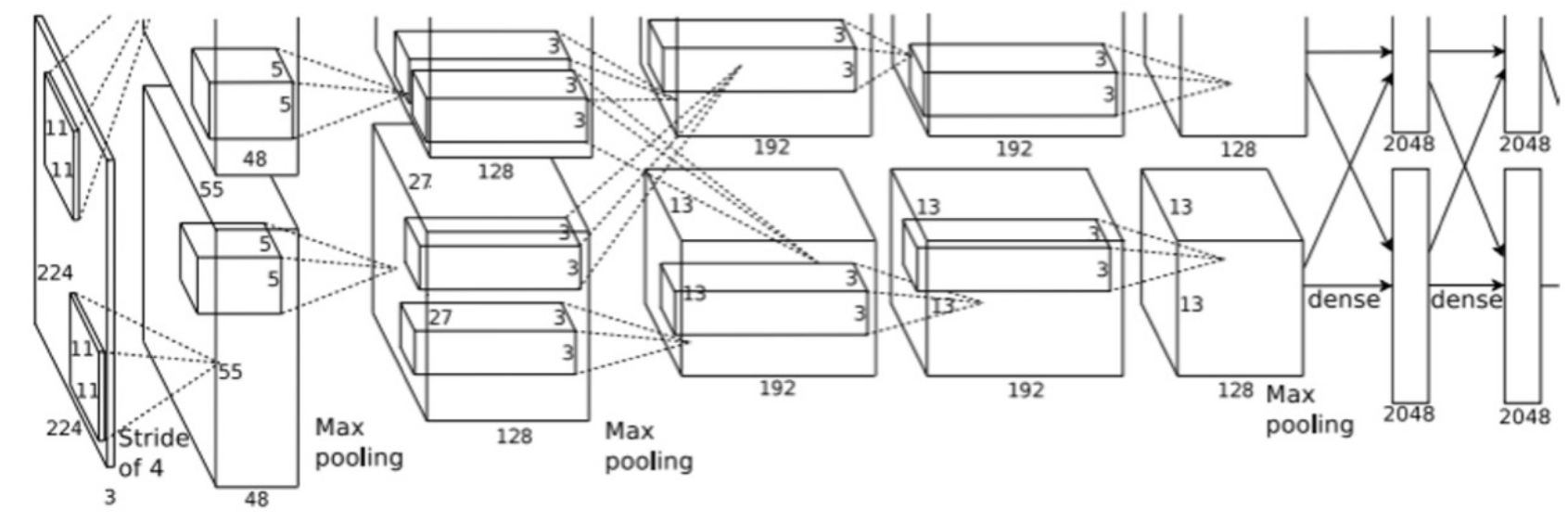
Chips: (x, y, w, h)
Chips: (x, y, w, h)
.....
Many numbers!

Need different numbers of output per image



Detecting Multiple Objects: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Hershey's: **No**

Flipz: **No**

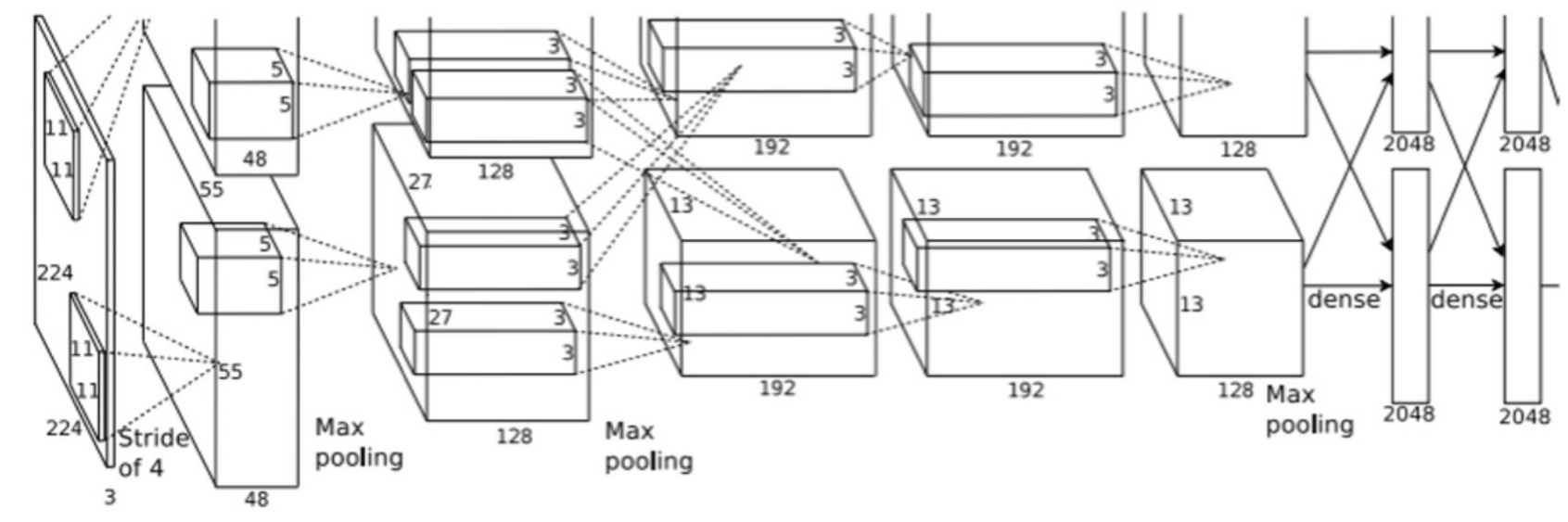
Reese's: **No**

Background: **Yes**



Detecting Multiple Objects: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Hershey's: **No**

Flipz: **Yes**

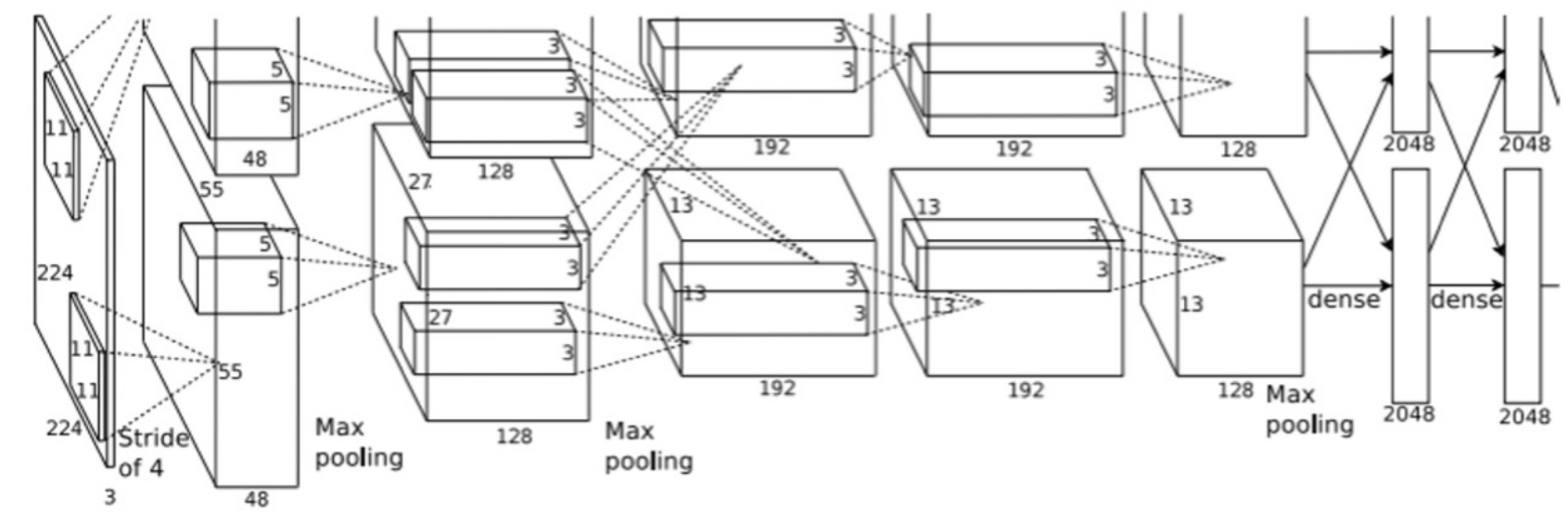
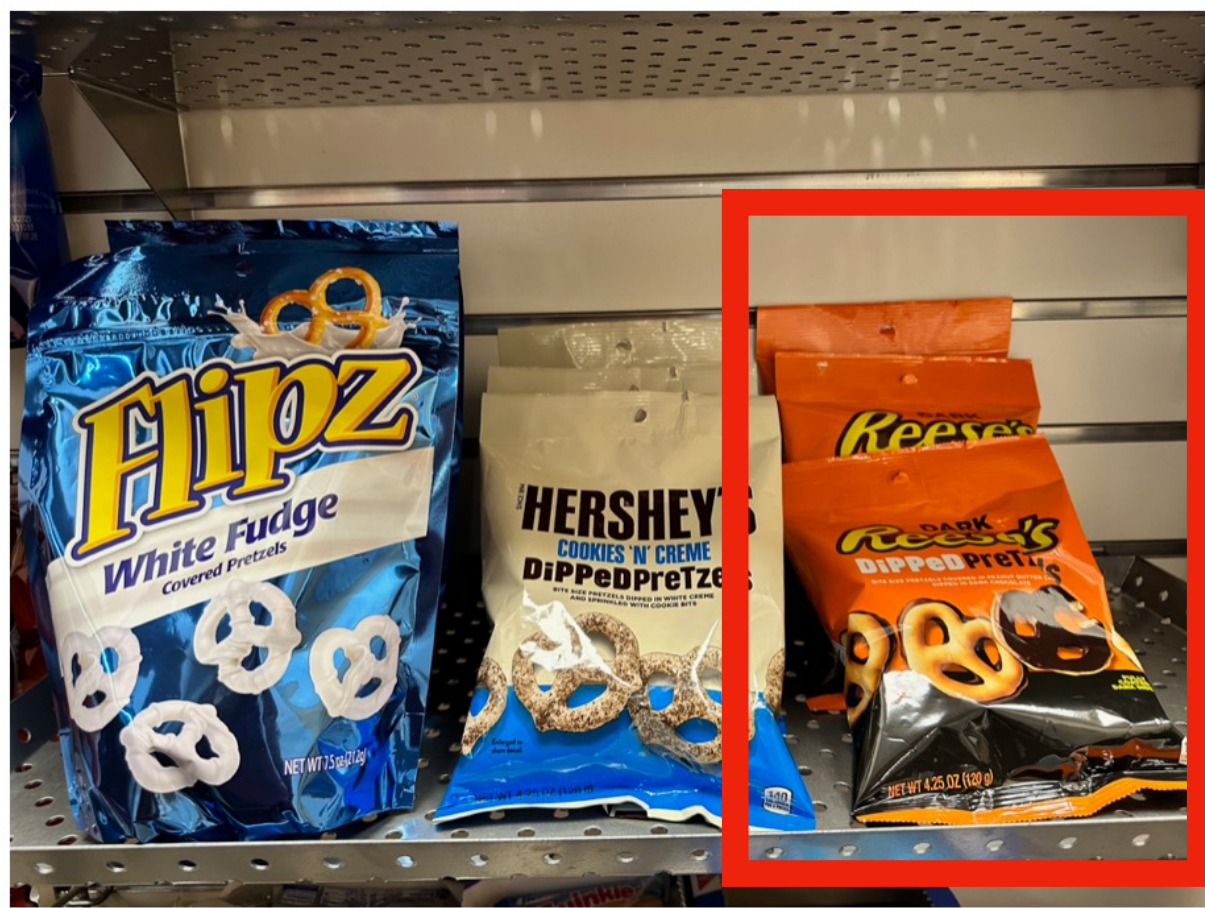
Reese's: **No**

Background: **No**



Detecting Multiple Objects: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Hershey's: **No**

Flipz: **No**

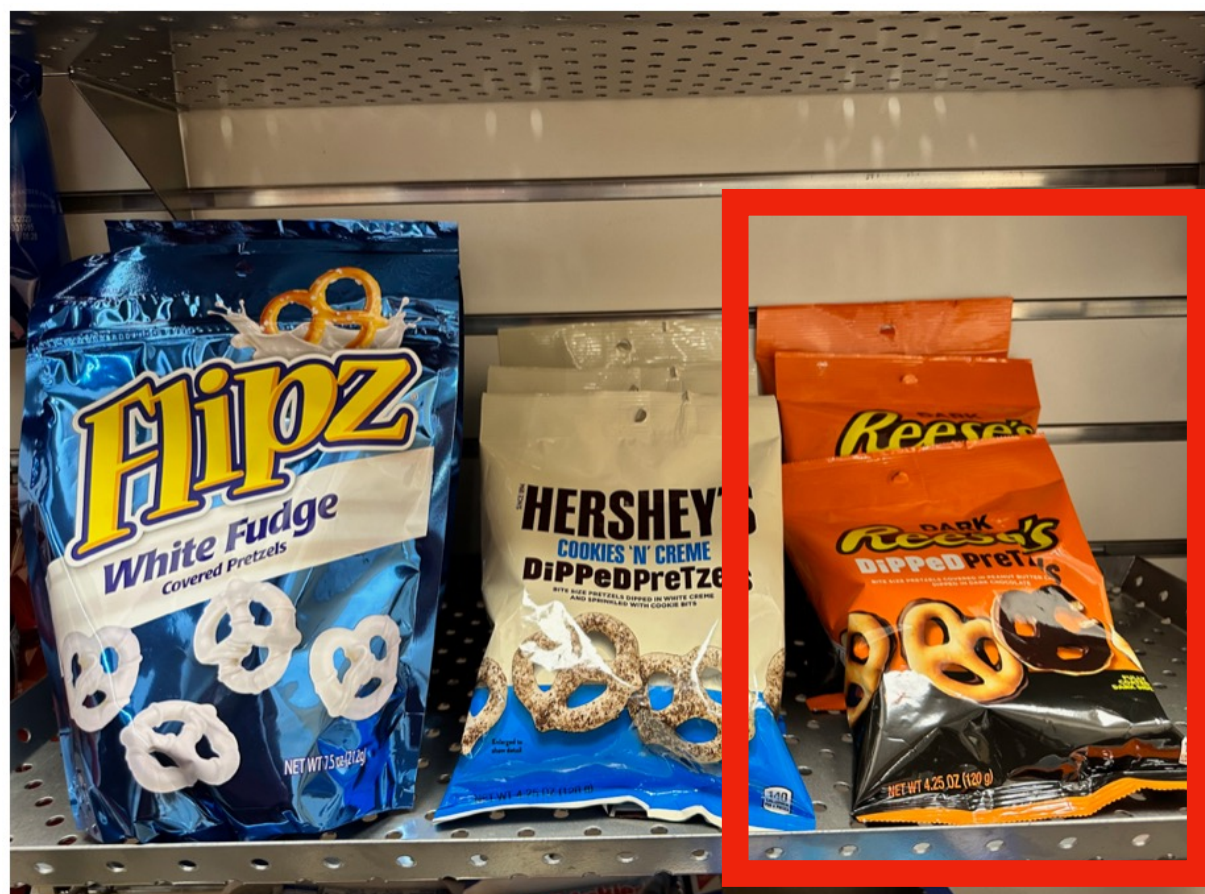
Reese's: **Yes**

Background: **No**



Detecting Multiple Objects: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background



Question: How many possible boxes are there in an image of size H x W?

Consider box of size h x w:
Possible x positions: $W - w + 1$
Possible y positions: $H - h + 1$
Possible positions:
 $(W-w+1) \times (H-h+1)$

Total possible boxes:

$$\sum_{h=1}^H \sum_{w=1}^W (W - w + 1)(H - h + 1)$$

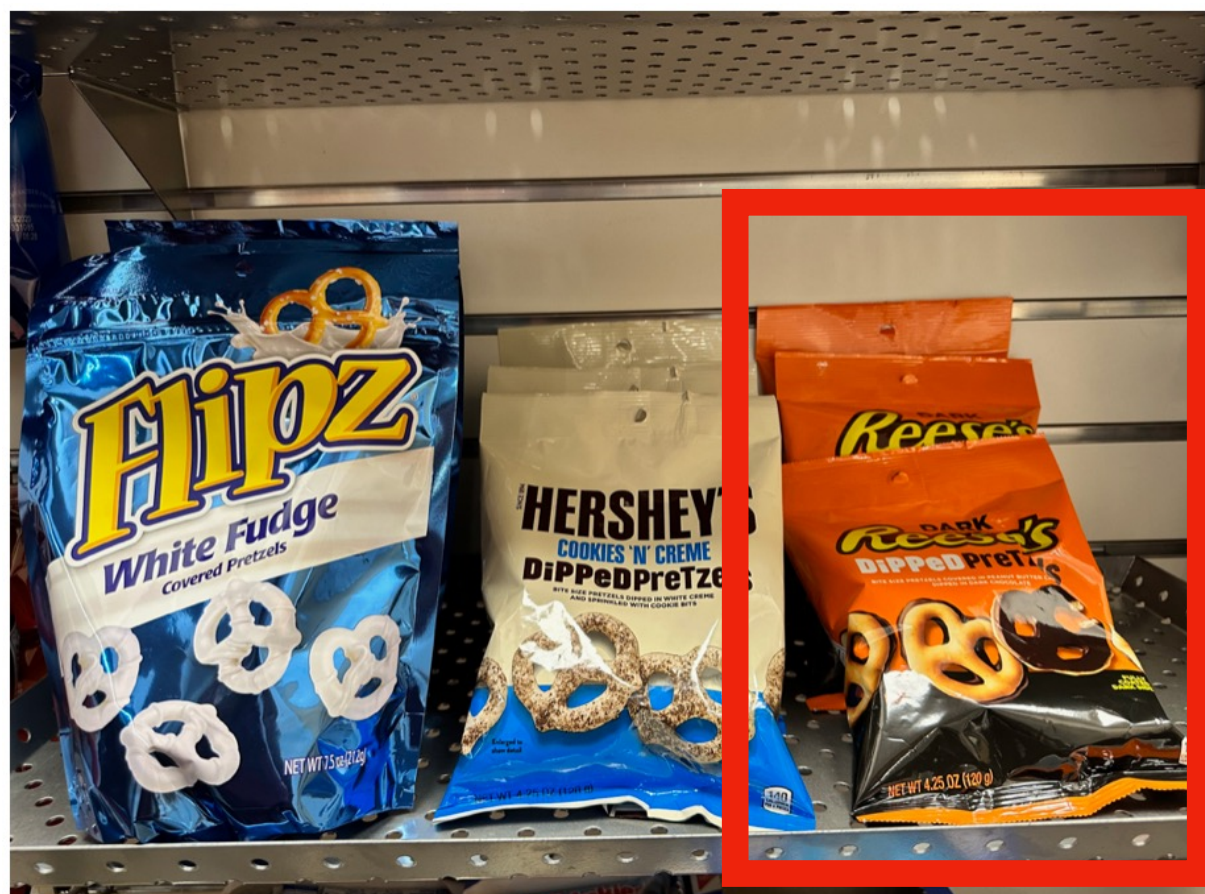
$$= \frac{H(H + 1)}{2} \frac{W(W + 1)}{2}$$



Detecting Multiple Objects: Sliding Window

Apply a CNN to many different crops of the image, CNN classifies each crop as object or background

800 x 600 image has ~58M boxes. No way we can evaluate them all



Question: How many possible boxes are there in an image of size H x W?

Consider box of size h x w:
Possible x positions: $W - w + 1$
Possible y positions: $H - h + 1$
Possible positions: $(W-w+1) \times (H-h+1)$

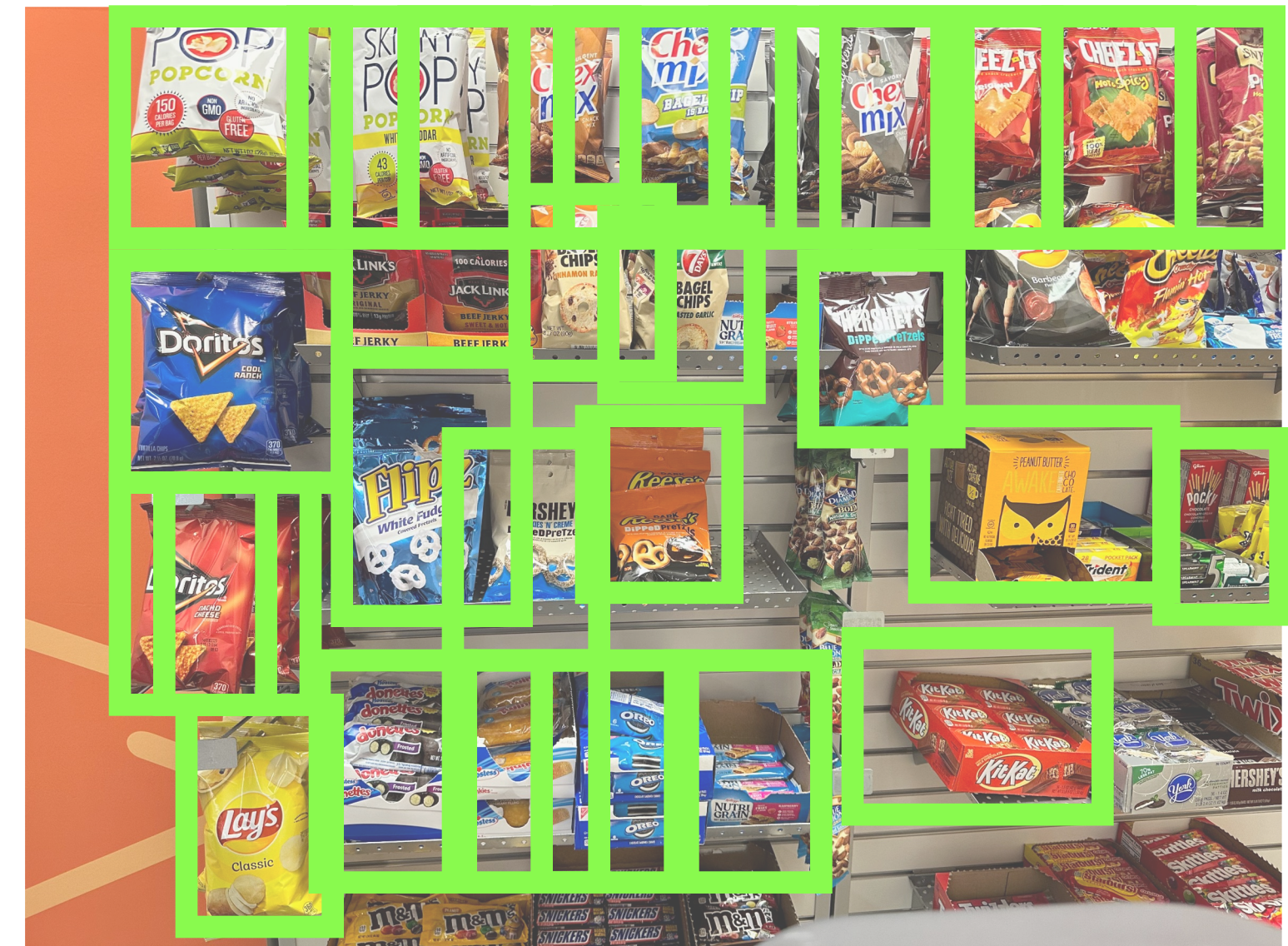
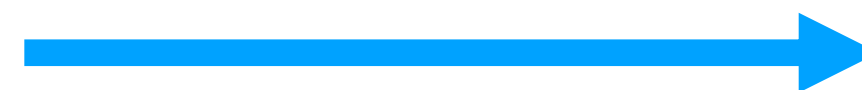
Total possible boxes:
$$\sum_{h=1}^H \sum_{w=1}^W (W - w + 1)(H - h + 1)$$

$$= \frac{H(H + 1)}{2} \frac{W(W + 1)}{2}$$



Region Proposals

- Find **a small set** of boxes that are likely to cover all objects
- Often based on **heuristics**: e.g. look for “blob-like” image regions
- **Relatively fast to run**; e.g. Selective Search gives 2000 region proposals in a few seconds on CPU





R-CNN: Region-Based CNN

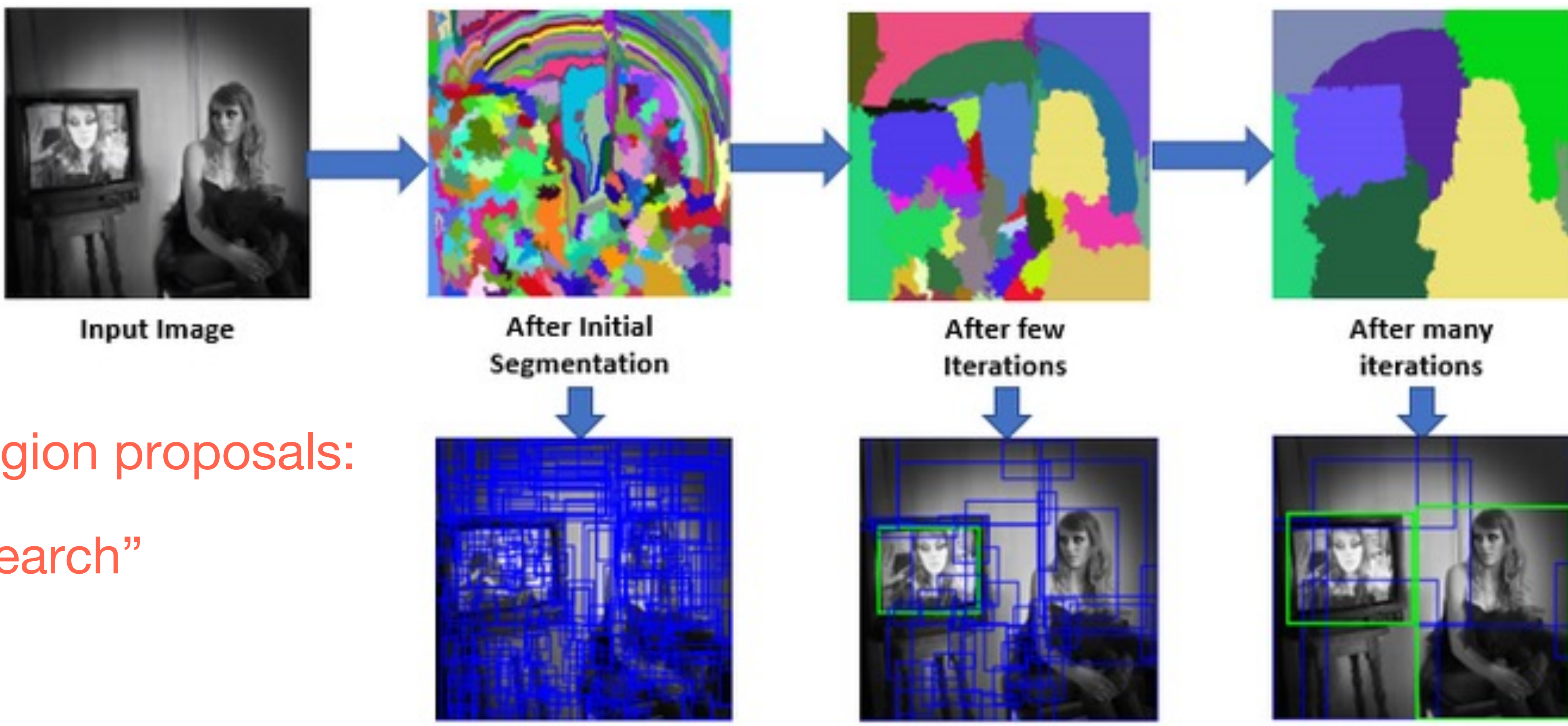
R-CNN: Region-Based CNN

Input
image





R-CNN: Region-Based CNN





R-CNN: Region-Based CNN

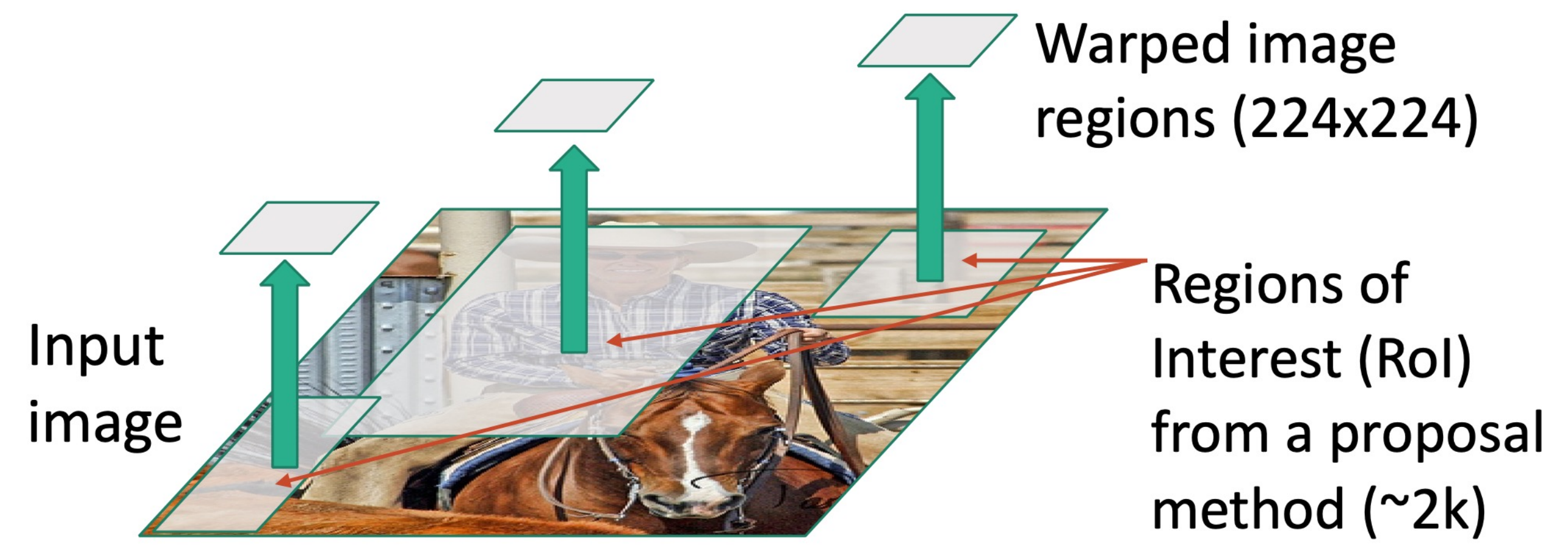
R-CNN: Region-Based CNN





R-CNN: Region-Based CNN

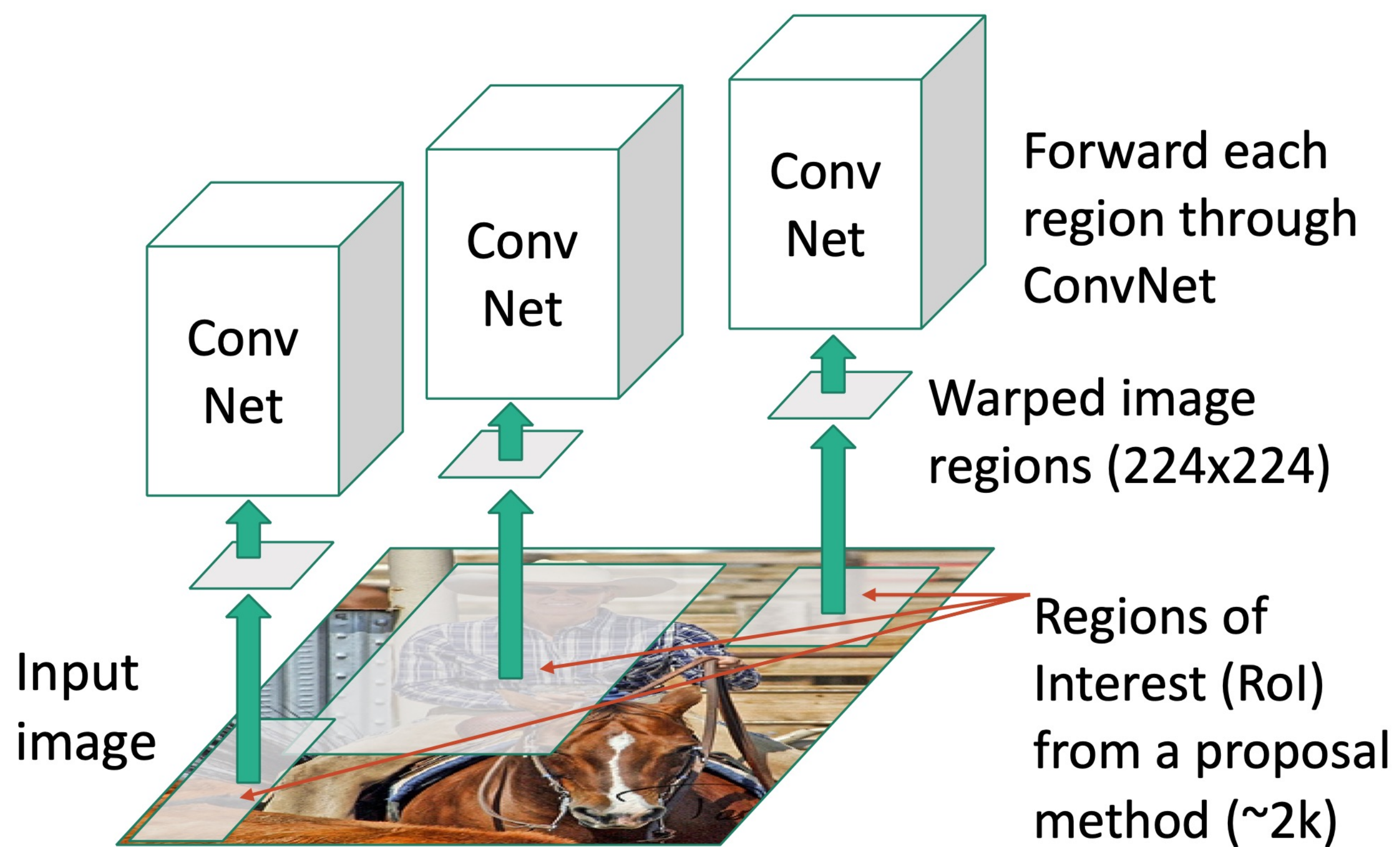
R-CNN: Region-Based CNN





R-CNN: Region-Based CNN

R-CNN: Region-Based CNN

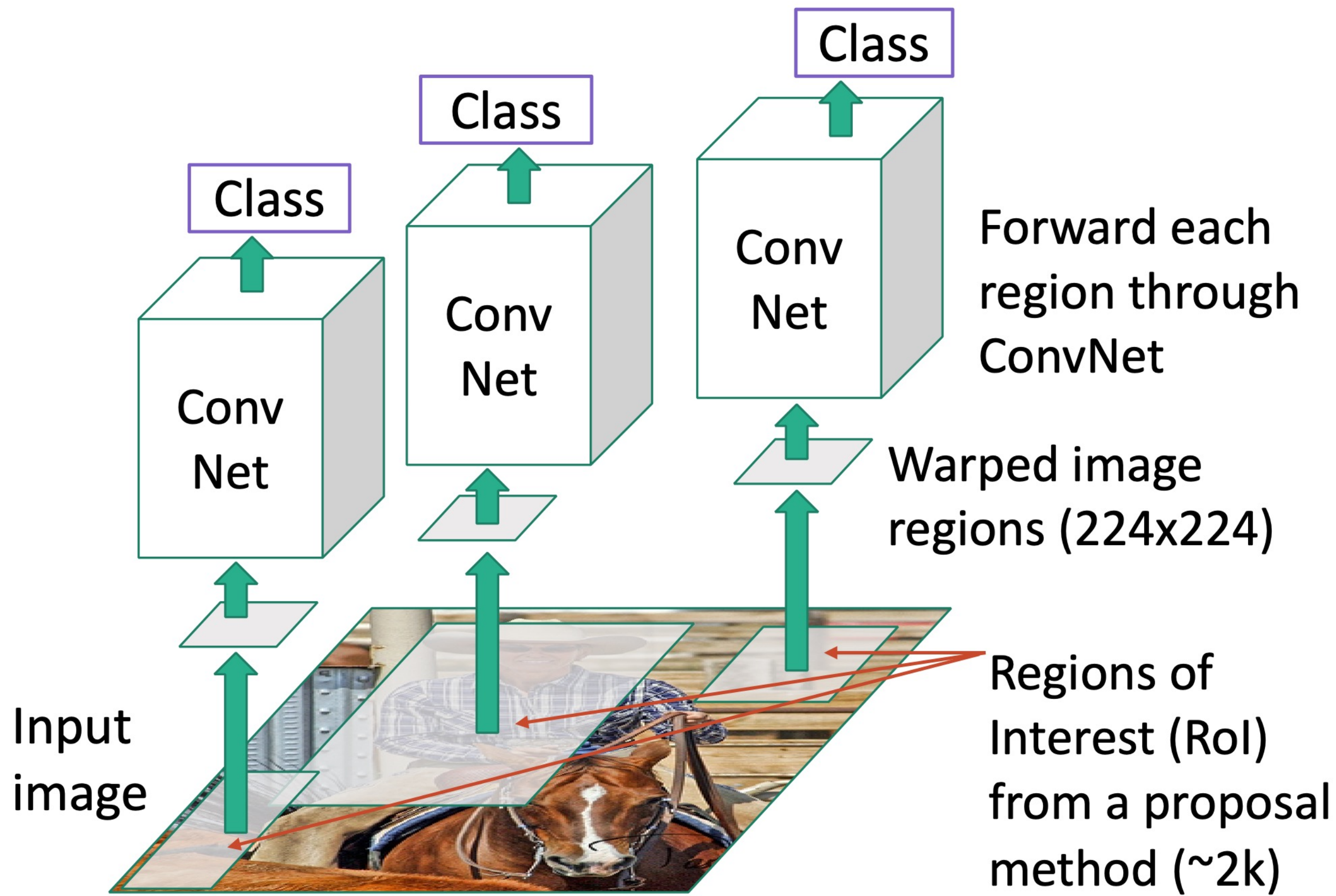




R-CNN: Region-Based CNN

R-CNN: Region-Based CNN

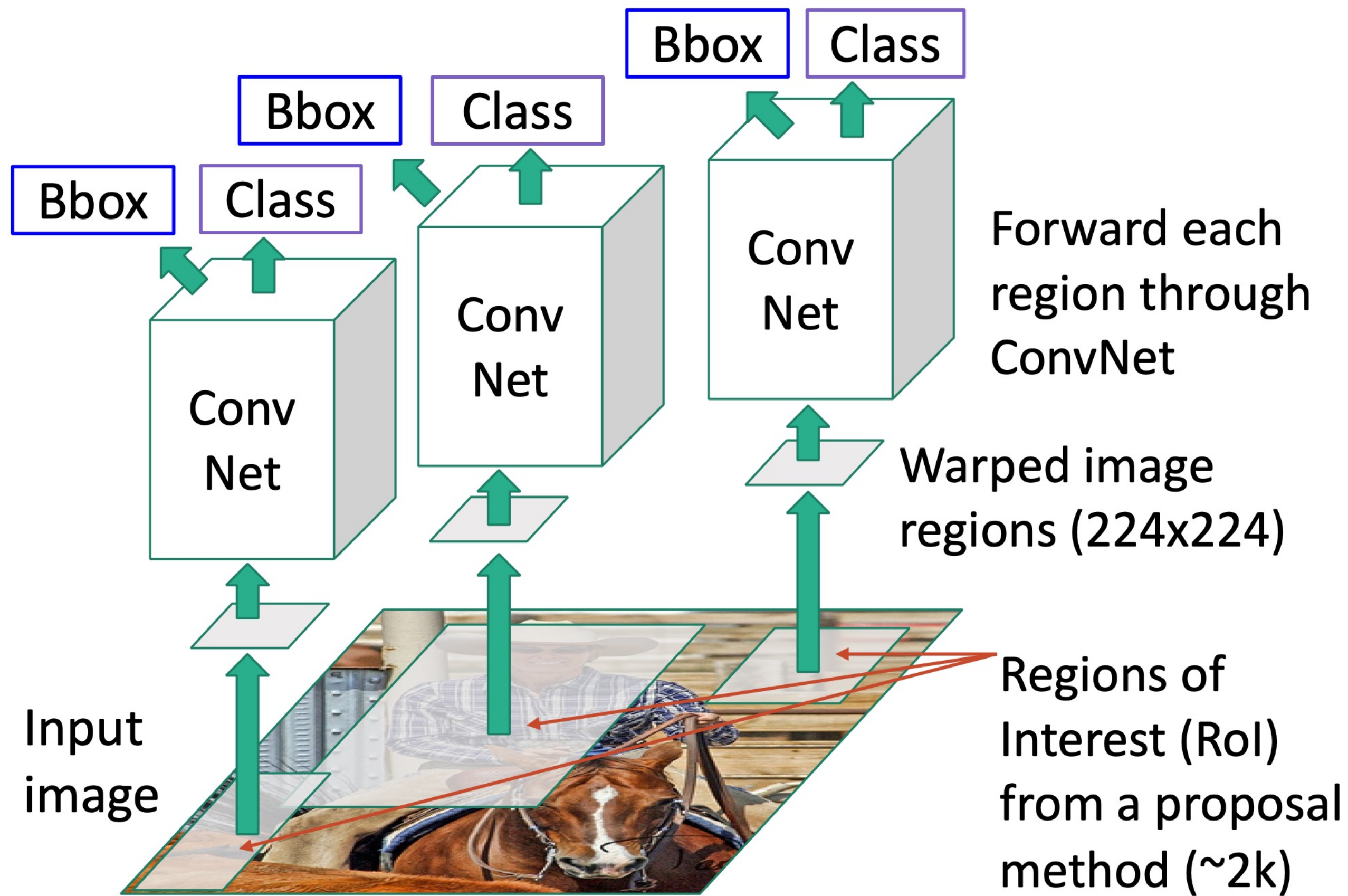
Classify each region





R-CNN: Region-Based CNN

R-CNN: Region-Based CNN



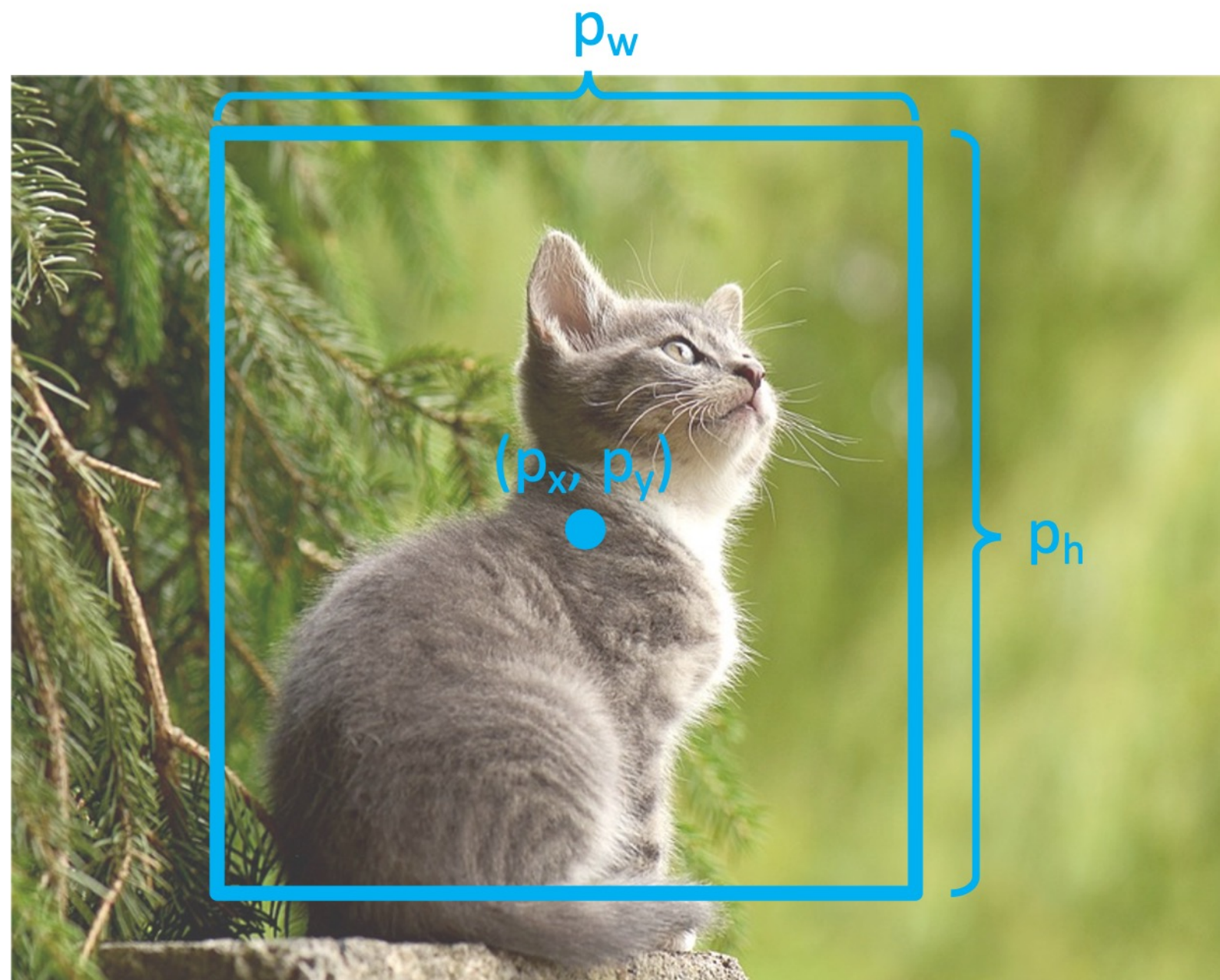
Classify each region

Bounding box regression:
Predict “transform” to correct the RoI: 4 numbers (t_x, t_y, t_h, t_w)



R-CNN: Box Regression

Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h



Model predicts a **transform** (t_x, t_y, t_w, t_h) to correct the region proposal



R-CNN: Box Regression

Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts a **transform** (t_x, t_y, t_w, t_h) to correct the region proposal

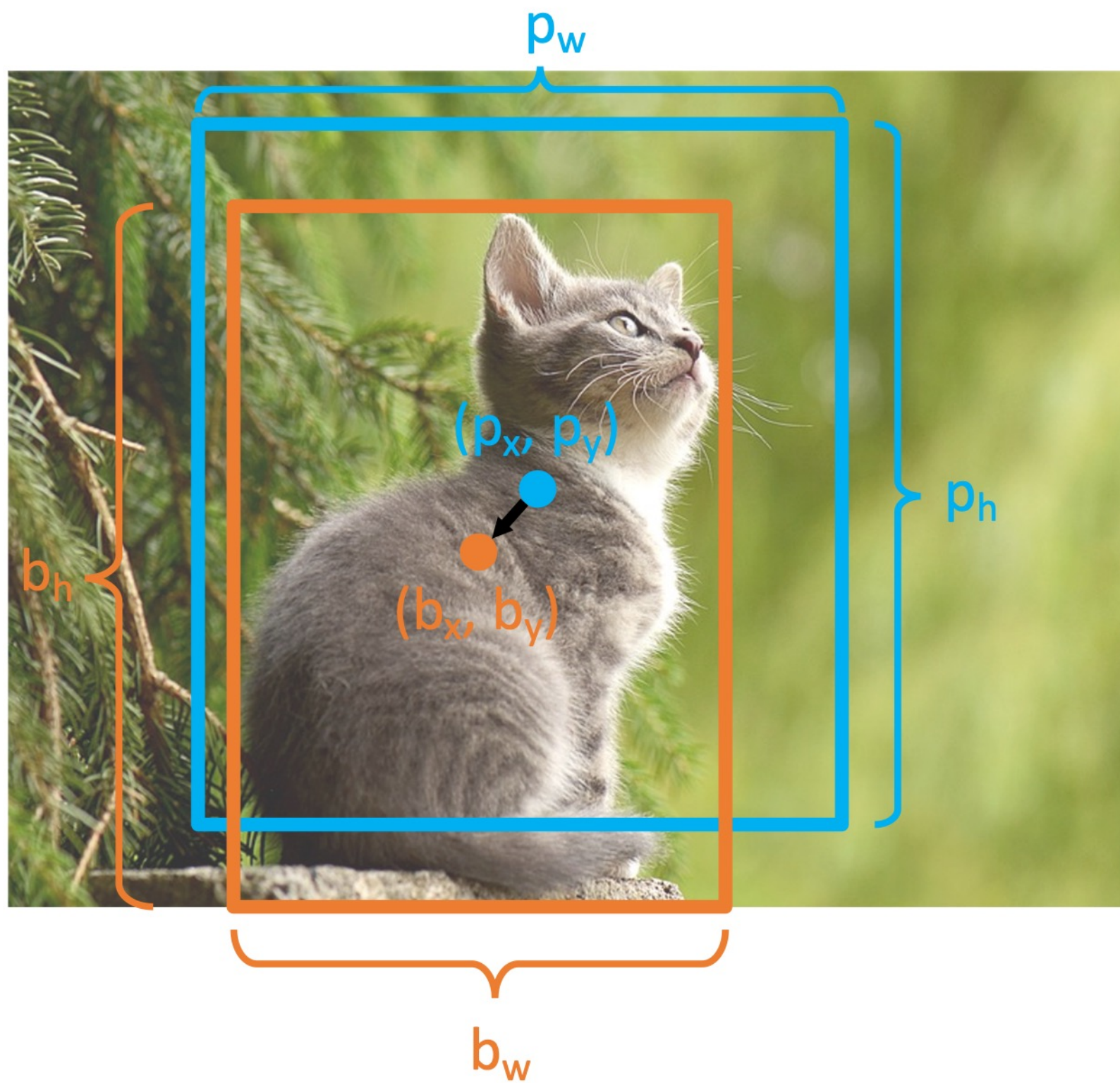
The **output box** is defined by:

$$b_x = p_x + p_w t_x \quad \text{Shift center by amount relative to proposal size}$$

$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w) \quad \text{Scale proposal; exp ensures that scaling factor is } > 0$$

$$b_h = p_h \exp(t_h)$$





R-CNN: Box Regression

Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts a **transform** (t_x, t_y, t_w, t_h) to correct the region proposal

The **output box** is defined by:

$$b_x = p_x + p_w t_x$$

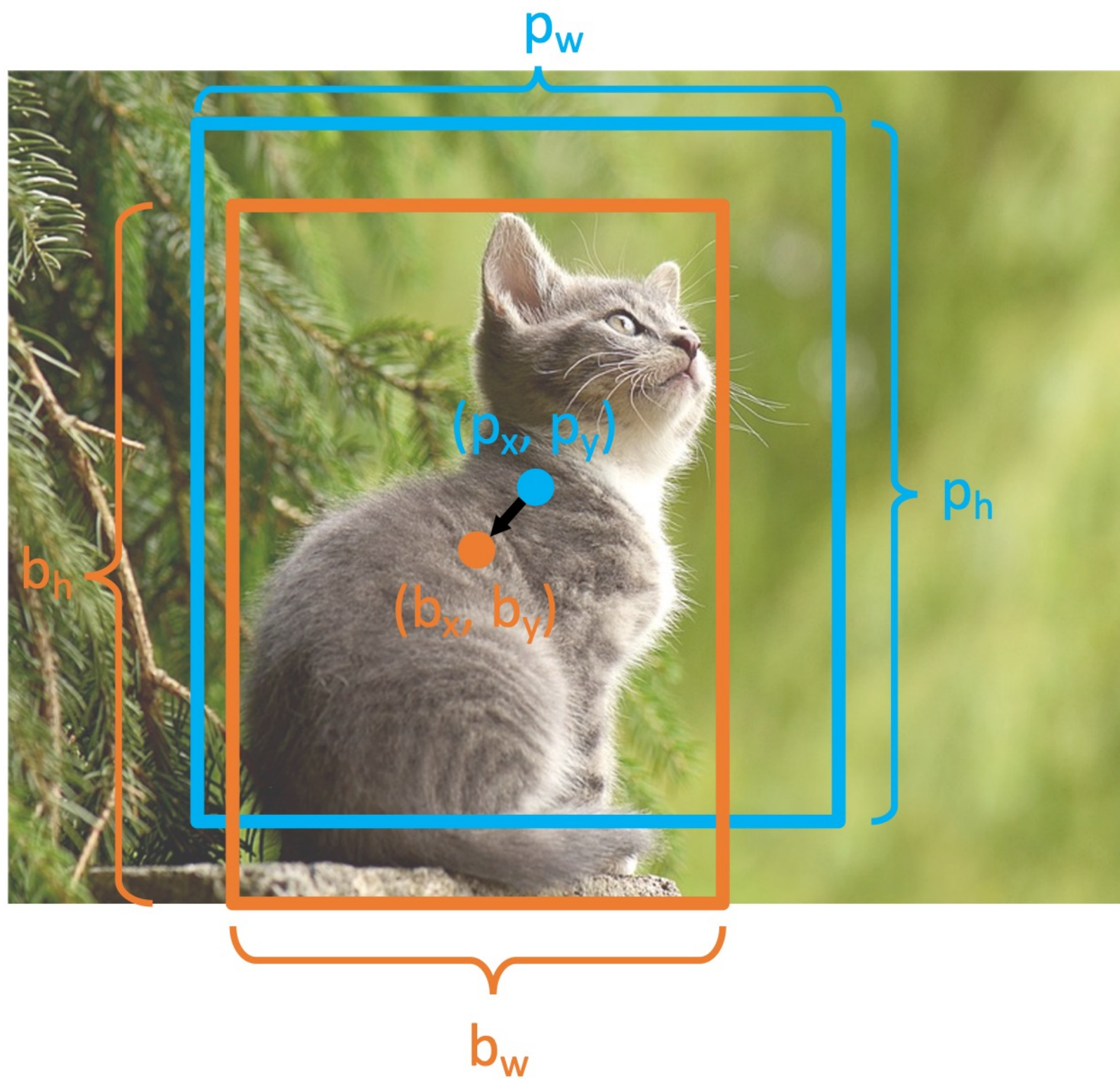
$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

When transform is 0, output = proposal

L2 regularization encourages leaving proposal unchanged





R-CNN: Box Regression

Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts a **transform** (t_x, t_y, t_w, t_h) to correct the region proposal

The **output box** is defined by:

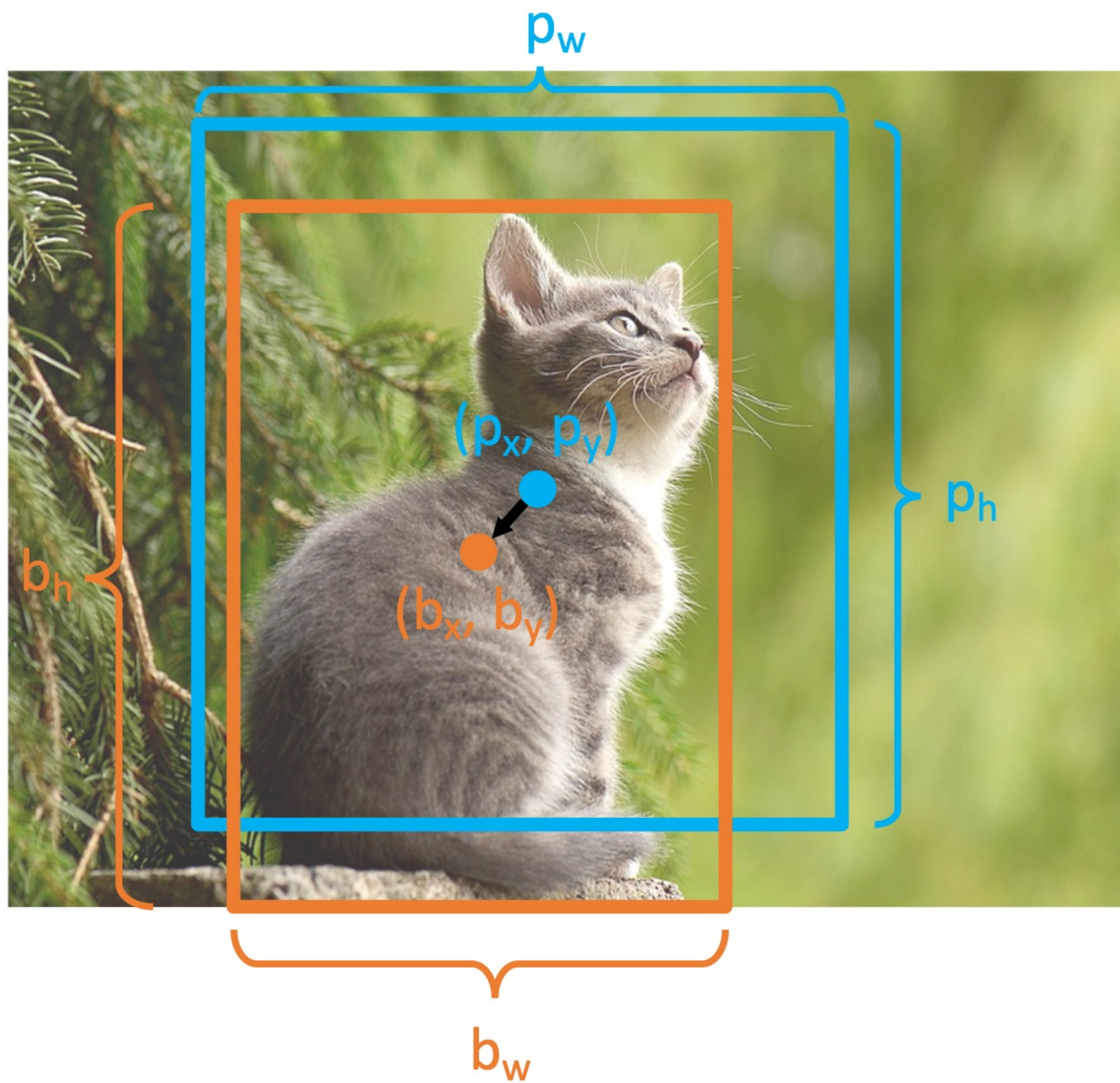
$$b_x = p_x + p_w t_x$$

$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

Scale / Translation invariance:
Transform encodes *relative* difference between proposal and output; important since CNN doesn't see absolute size or position after cropping

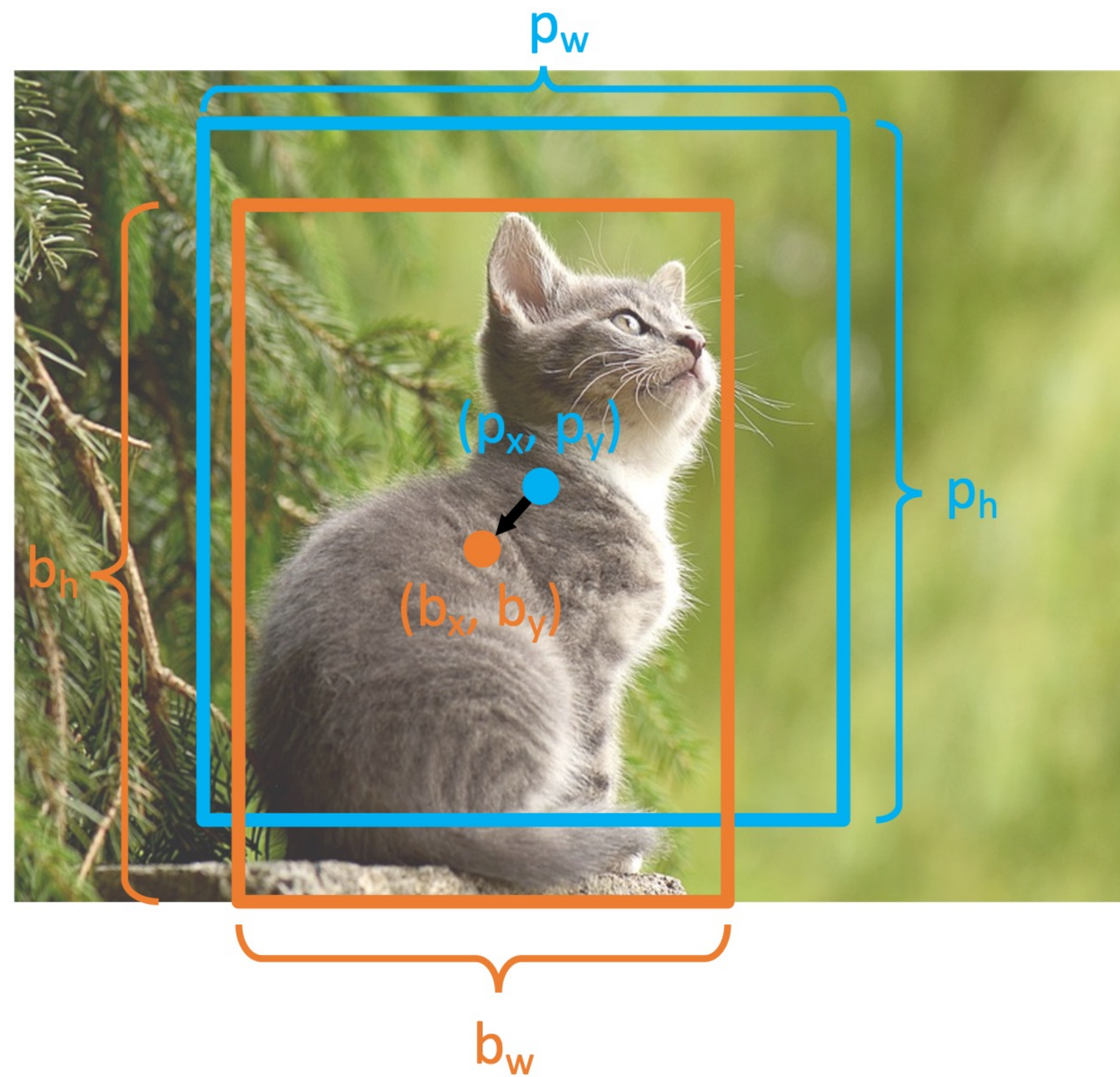




R-CNN: Box Regression

Consider a **region proposal** with center (p_x, p_y) , width p_w , height p_h

Model predicts a **transform** (t_x, t_y, t_w, t_h) to correct the region proposal



The **output box** is defined by:

$$b_x = p_x + p_w t_x$$

$$b_y = p_y + p_h t_y$$

$$b_w = p_w \exp(t_w)$$

$$b_h = p_h \exp(t_h)$$

Given **proposal** and **target output**, we can solve for the **transform** the network should output:

$$t_x = (b_x - p_x) / p_w$$

$$t_y = (b_y - p_y) / p_h$$

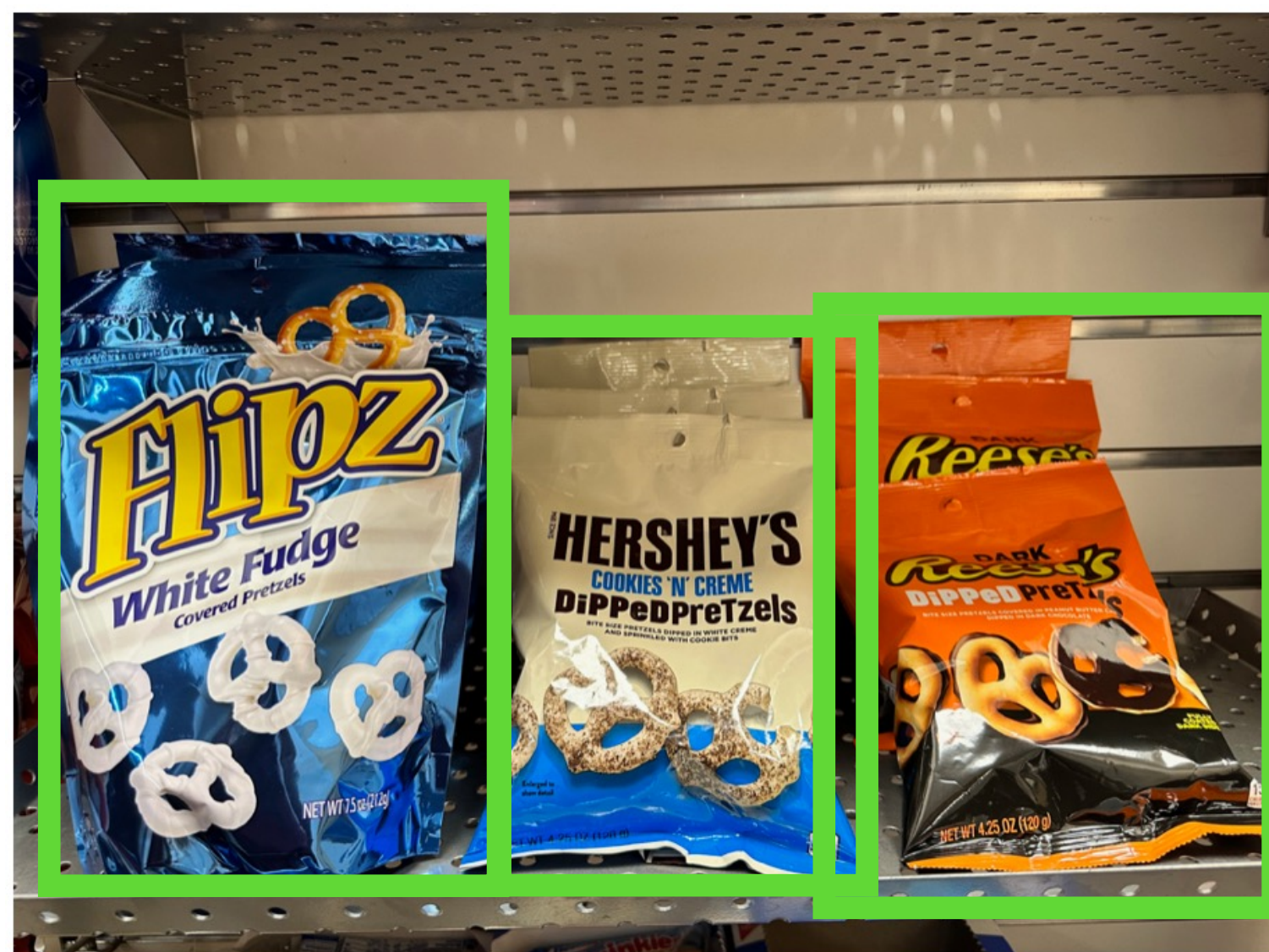
$$t_w = \log(b_w / p_w)$$

$$t_h = \log(b_h / p_h)$$



R-CNN: Training

Input Image



Ground Truth



R-CNN: Training

Input Image



Ground Truth

Region Proposals



R-CNN: Training

Input Image



Ground Truth	Positive
--------------	----------

Neutral	Negative
---------	----------



R-CNN: Training

Input Image



Ground Truth	Positive
Neutral	Negative

Categorize each region proposal as **positive**, **negative** or neutral based on overlap with the Ground truth boxes:

Positive: > 0.5 IoU with a GT box

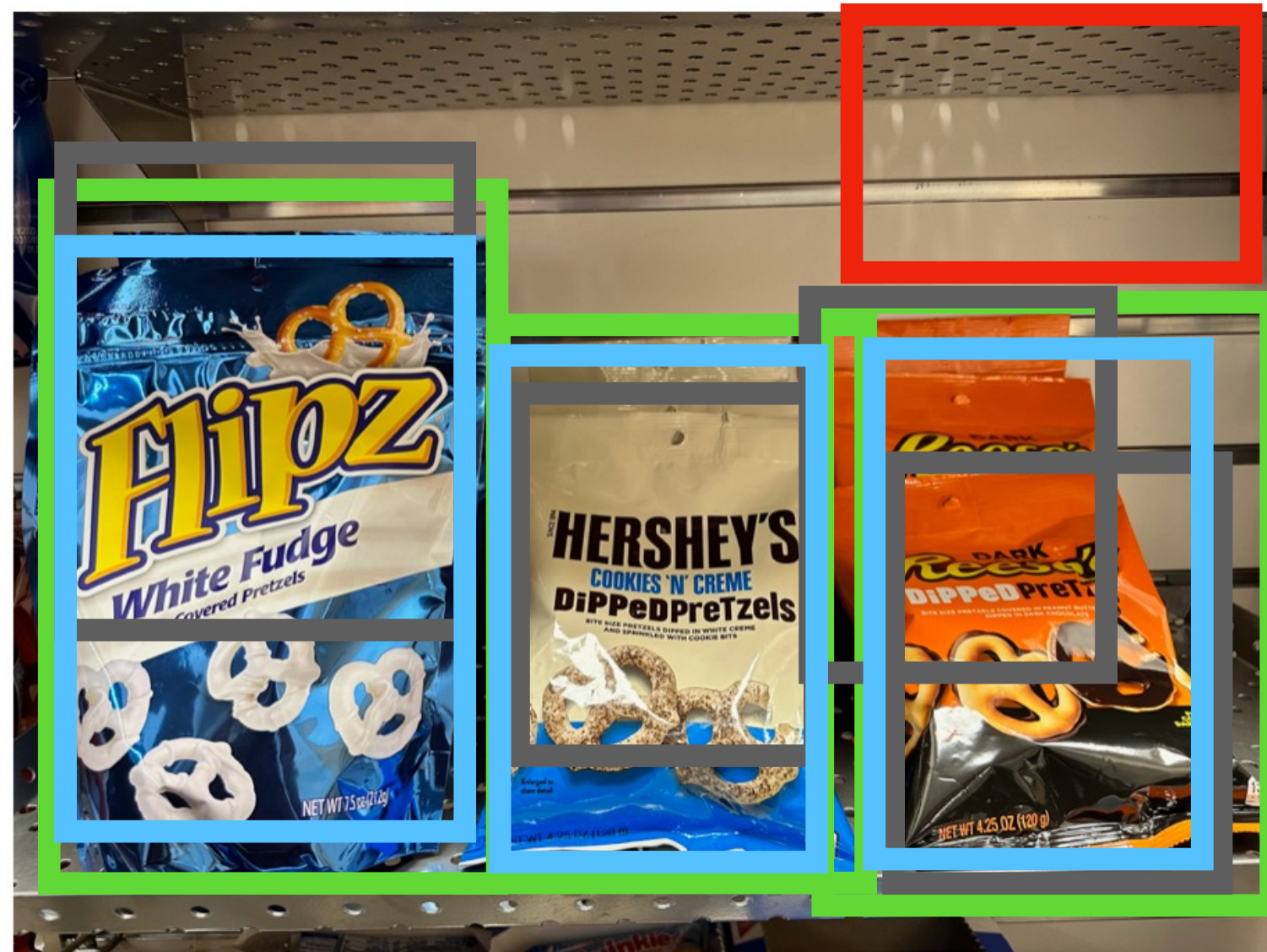
Negative: < 0.3 IoU with all GT boxes

Neutral: between 0.3 and 0.5 IoU with GT boxes



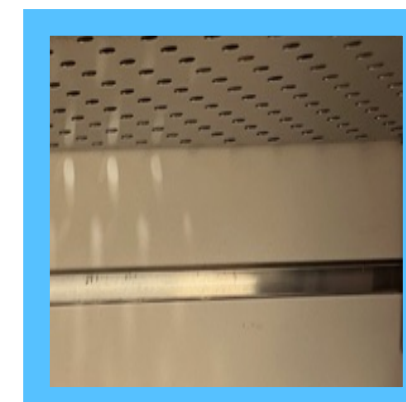
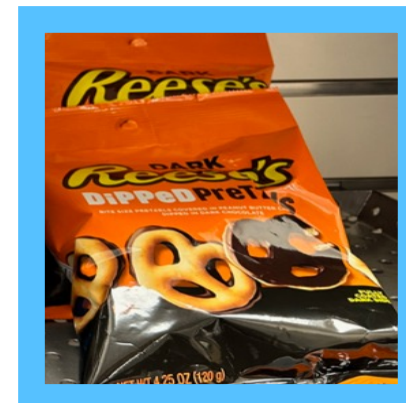
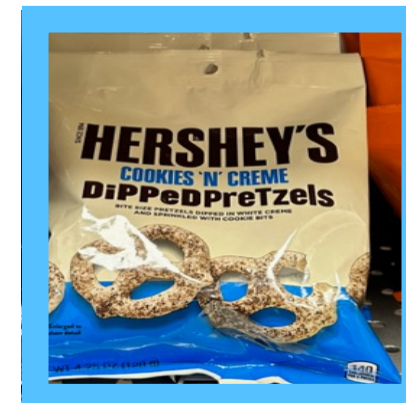
R-CNN: Training

Input Image



Ground Truth Positive

Neutral Negative



Crop pixels from each positive and negative proposal, resize to 224 x 224

Run each region through CNN

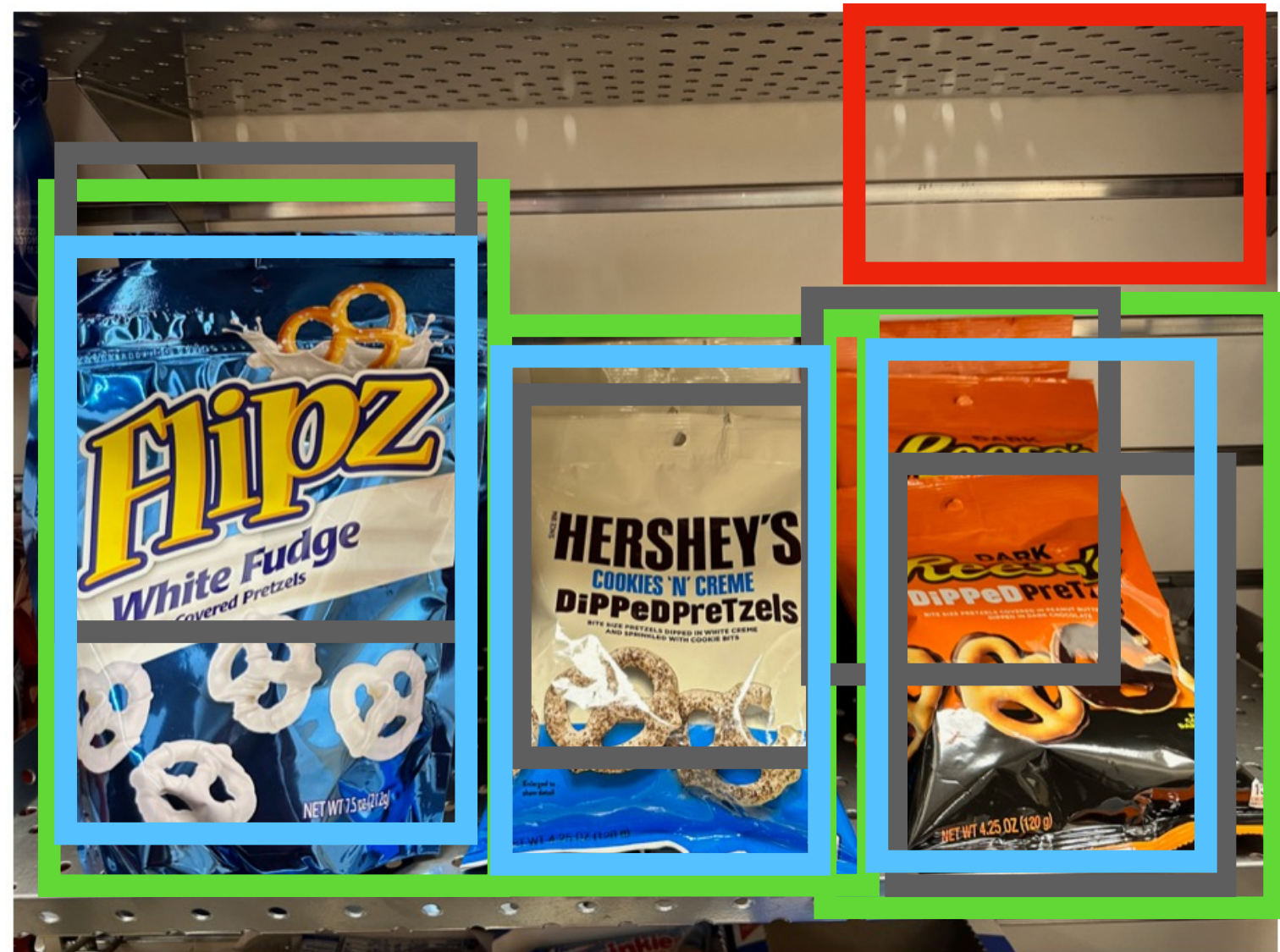
Positive regions: predict class and transform

Negative regions: just predict class



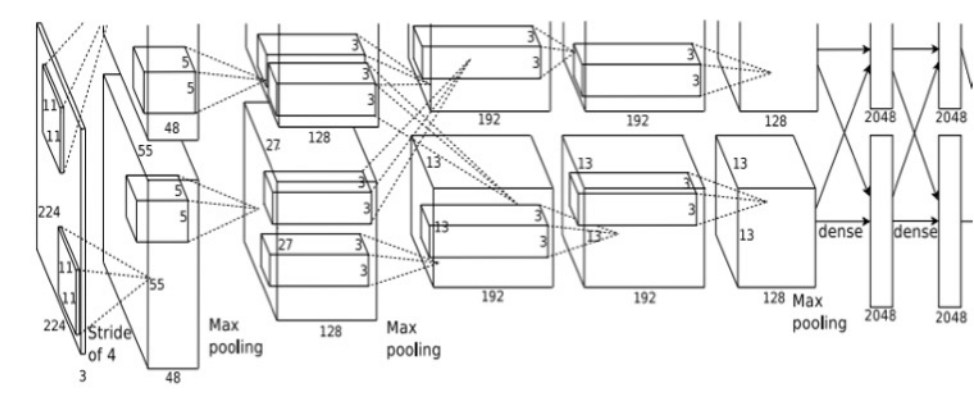
R-CNN: Training

Input Image



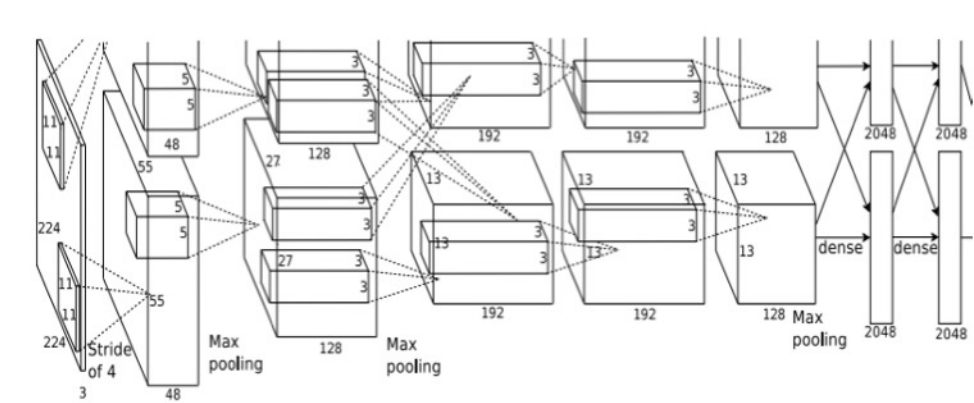
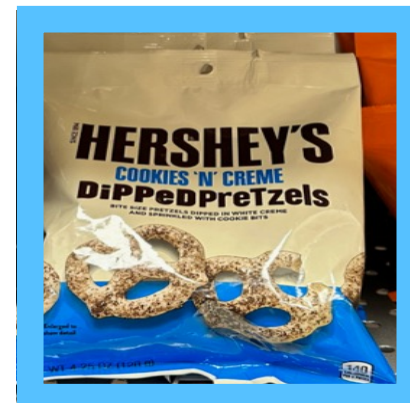
Ground Truth Positive

Neutral Negative



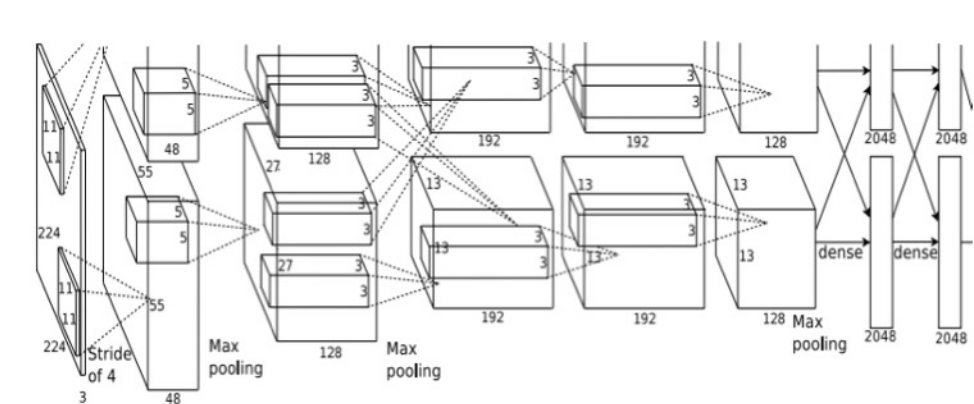
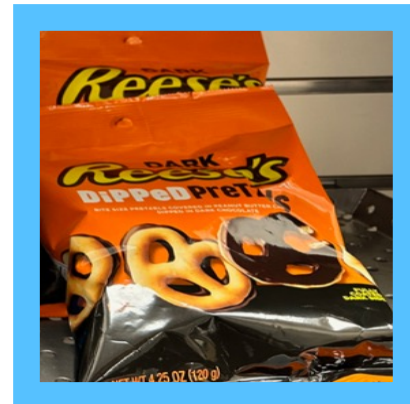
Class target: Flipz

Box target:



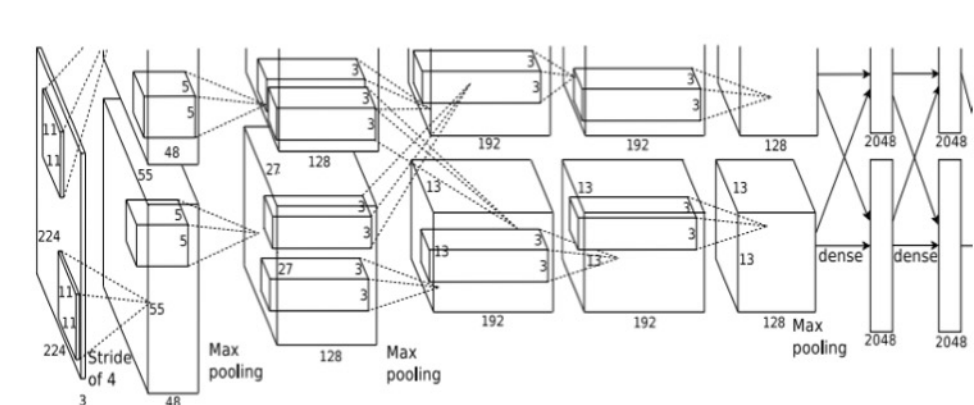
Class target: Hershey's

Box target:



Class target: Reese's

Box target:



Class target: Background

Box target: None

Run each region through CNN
Positive regions: predict class and transform
Negative regions: just predict class





R-CNN: Test time

Input Image



Region Proposals

Run proposal method:

1. Run CNN on each proposal to get class scores, transforms
2. Threshold class scores to get a set of detections

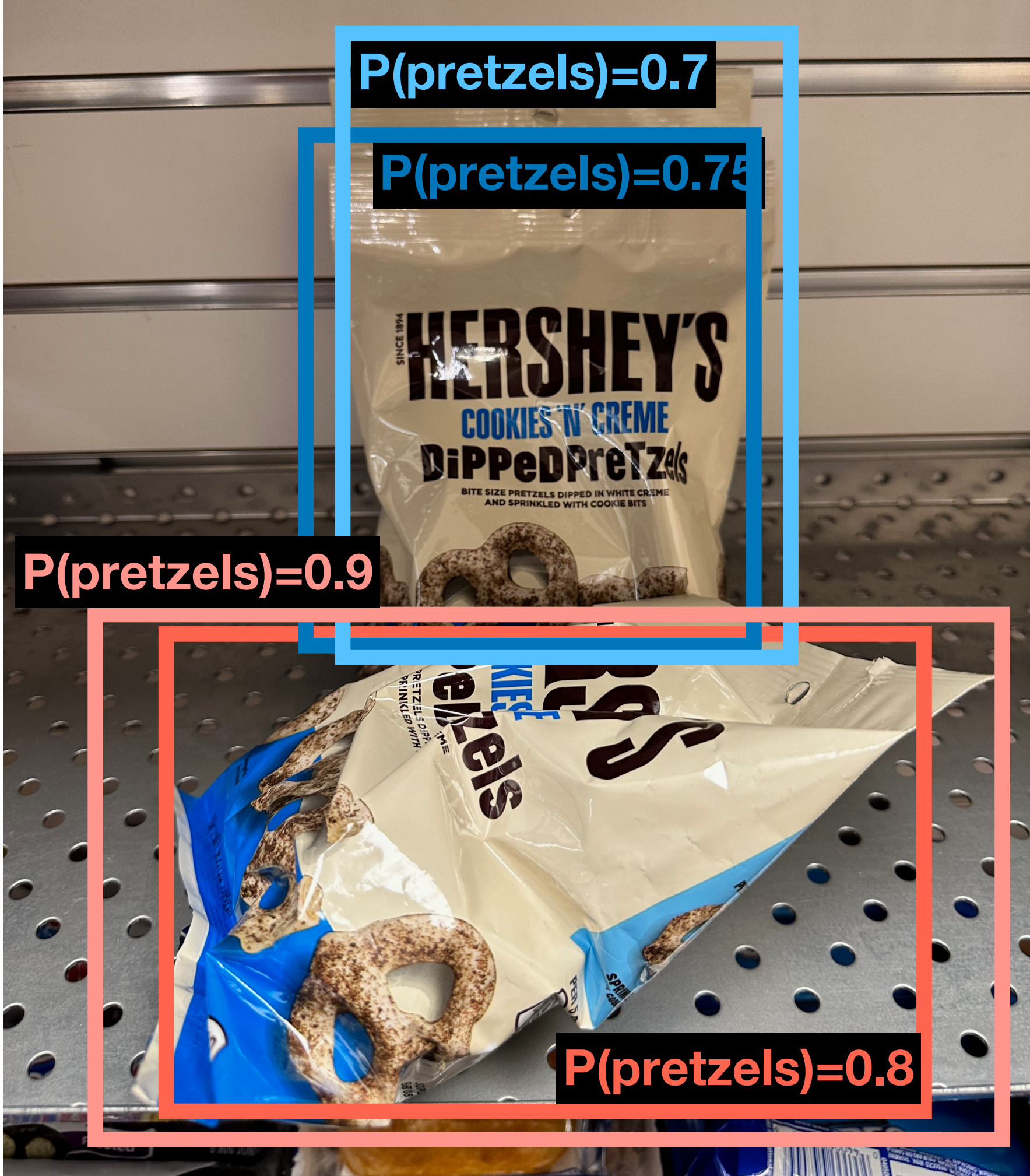
2 Problems:

1. CNN often outputs overlapping boxes
2. How to set thresholds?



Overlapping Boxes

Problem: Object detectors often output many overlapping detections



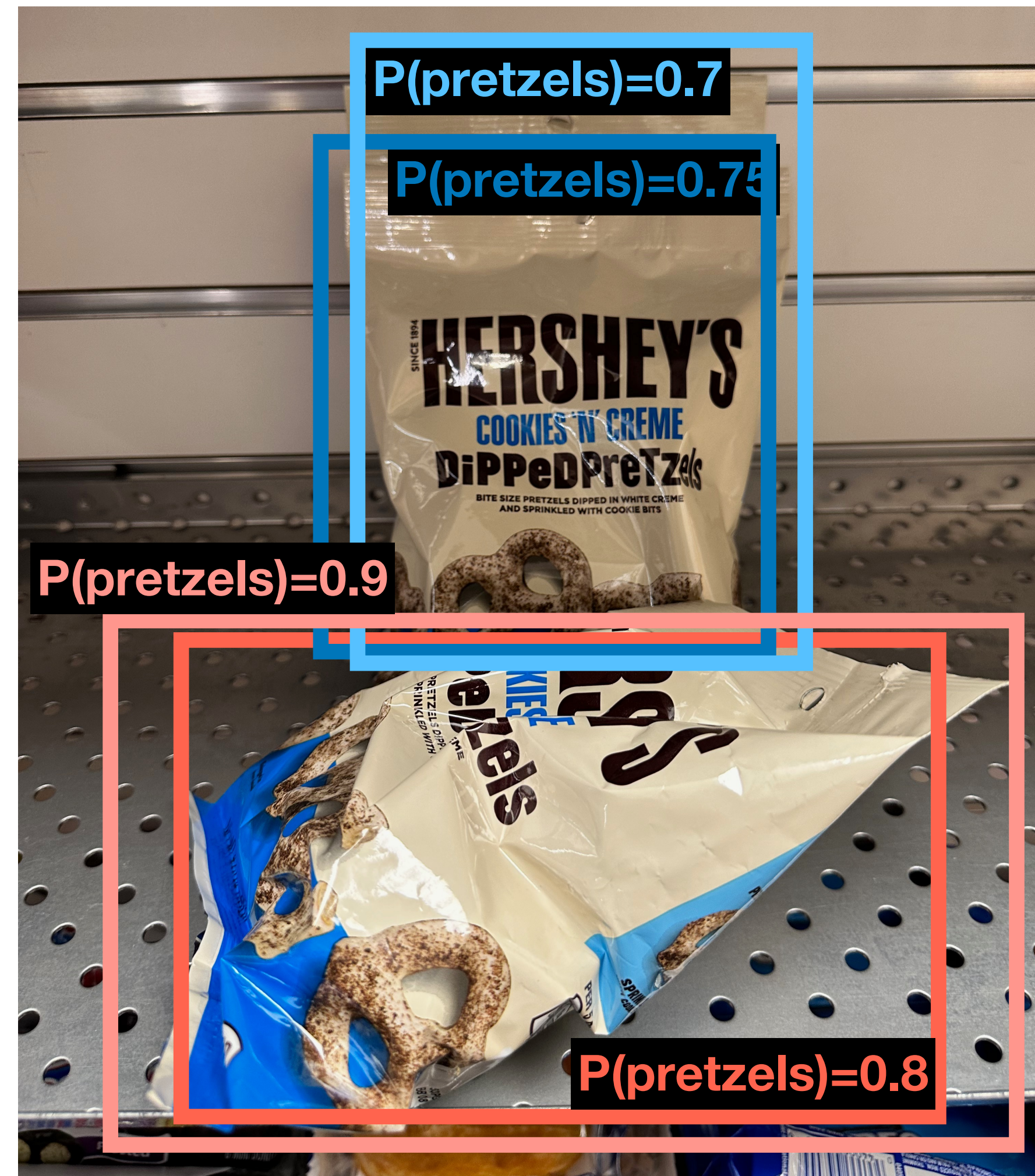


Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections

Solution: Post-process raw detections using Non-Max Suppression (NMS)

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $IoU > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1





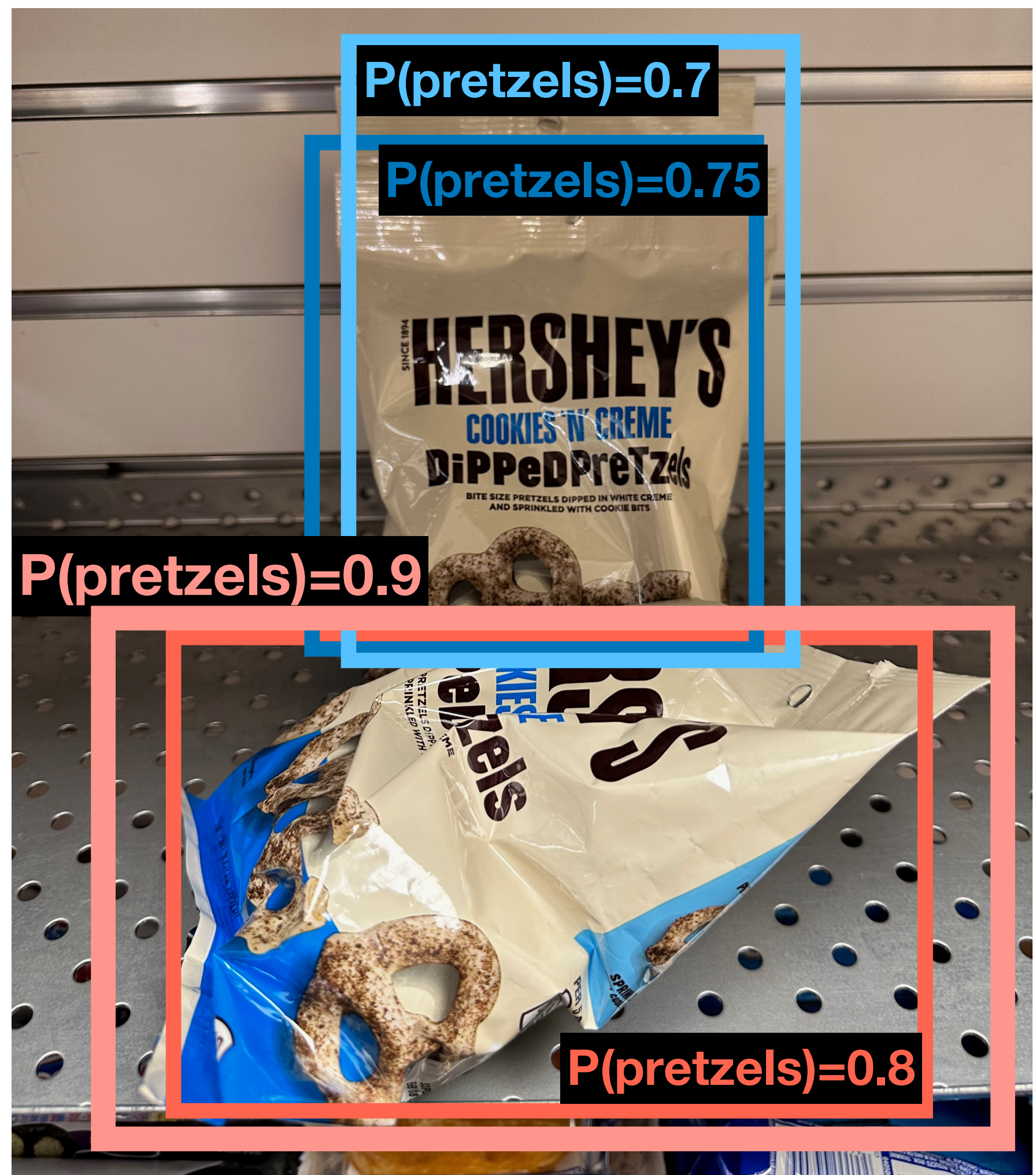
Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections

Solution: Post-process raw detections using Non-Max Suppression (NMS)

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $IoU > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1

$IoU(\text{red}, \text{red}) = 0.8$
 $IoU(\text{red}, \text{blue}) = 0.03$
 $IoU(\text{red}, \text{light blue}) = 0.05$





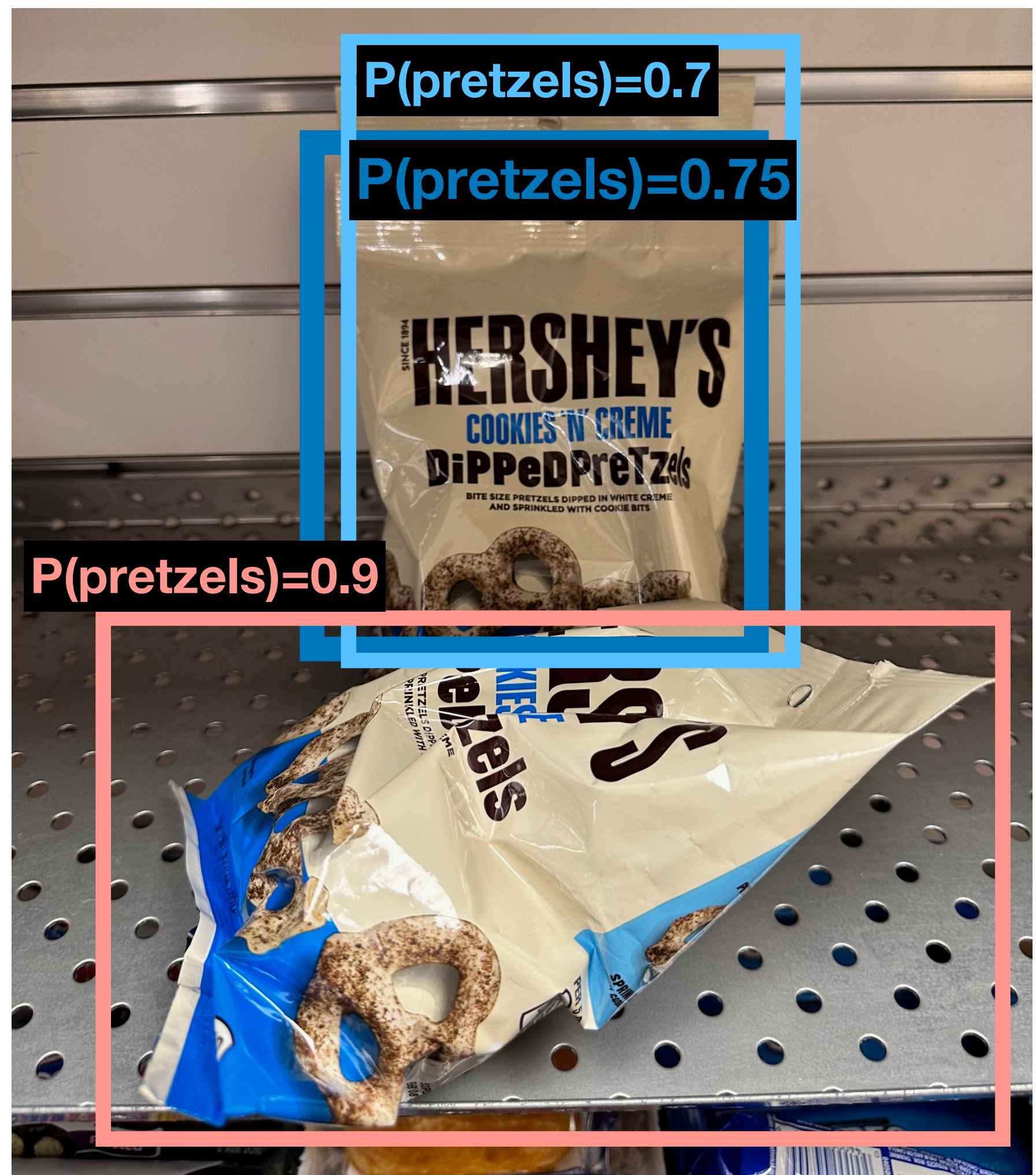
Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections

Solution: Post-process raw detections using Non-Max Suppression (NMS)

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $IoU > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1

$$IoU(\text{dark blue box}, \text{light blue box}) = 0.85$$



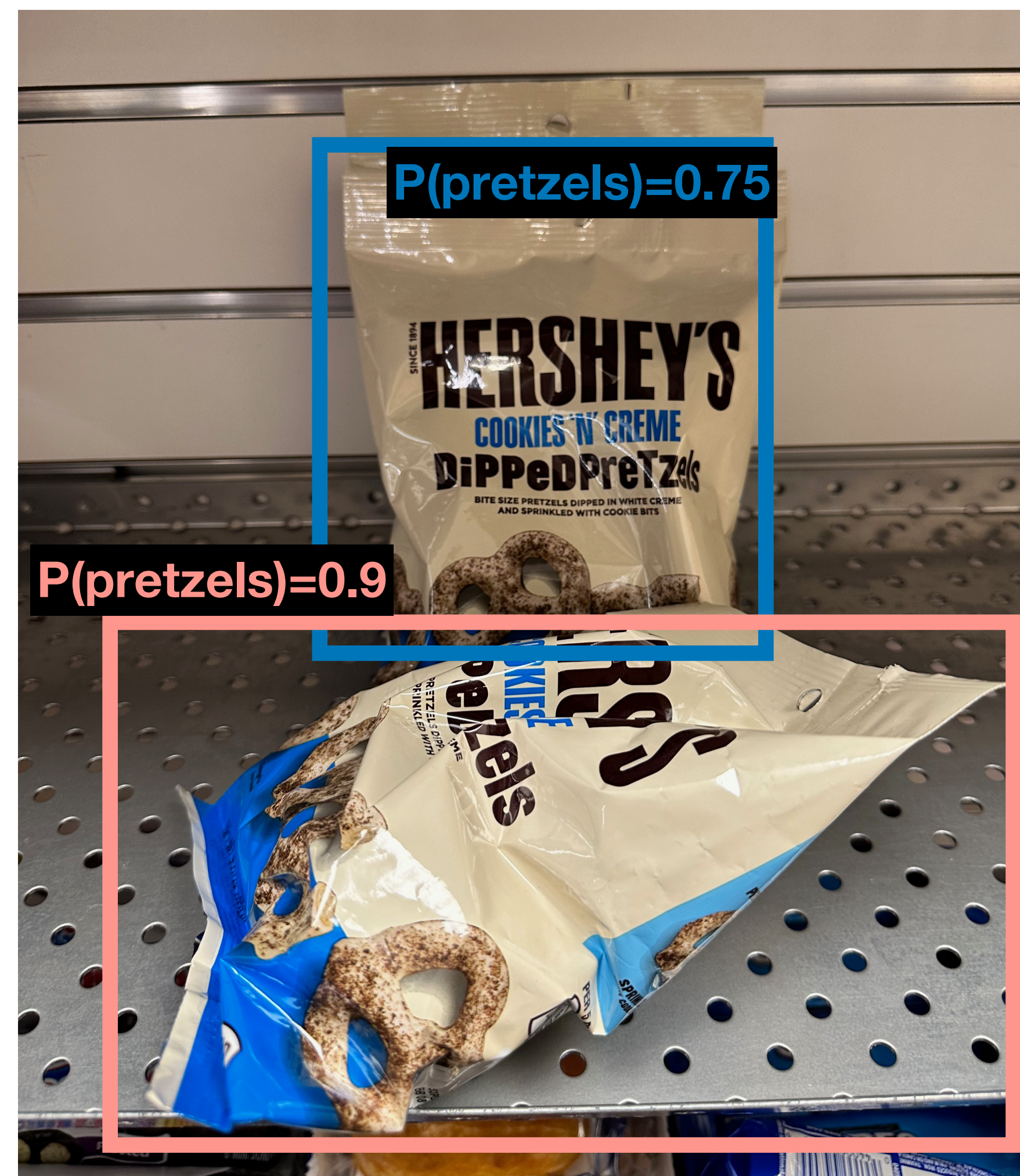


Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections

Solution: Post-process raw detections using Non-Max Suppression (NMS)

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $IoU > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1





Overlapping Boxes: Non-Max Suppression (NMS)

Problem: Object detectors often output many overlapping detections

Solution: Post-process raw detections using Non-Max Suppression (NMS)

1. Select next highest-scoring box
2. Eliminate lower-scoring boxes with $IoU > \text{threshold}$ (e.g. 0.7)
3. If any boxes remain, GOTO 1

Problem: NMS may eliminate “good” boxes when objects are highly overlapping... no good solution





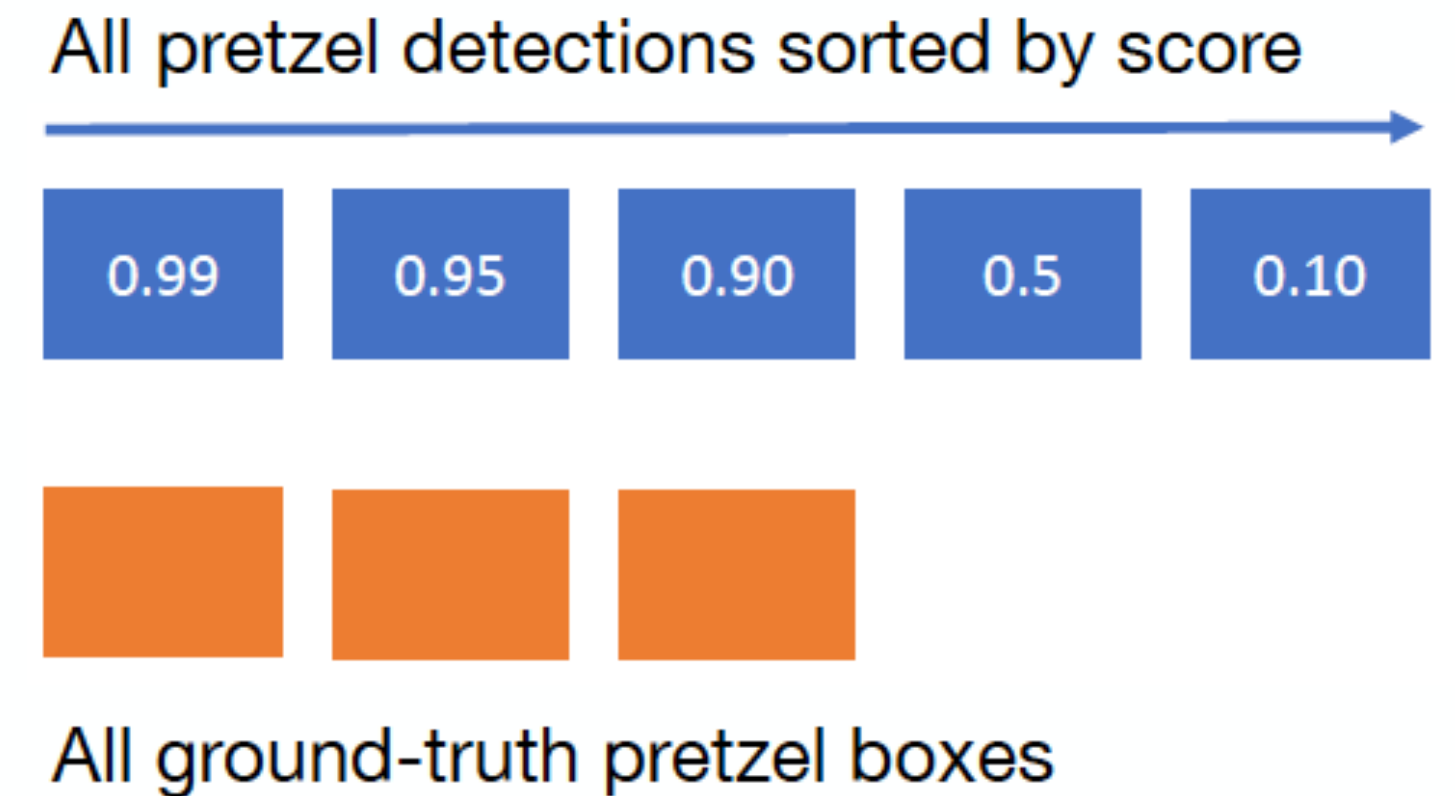
Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve



Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

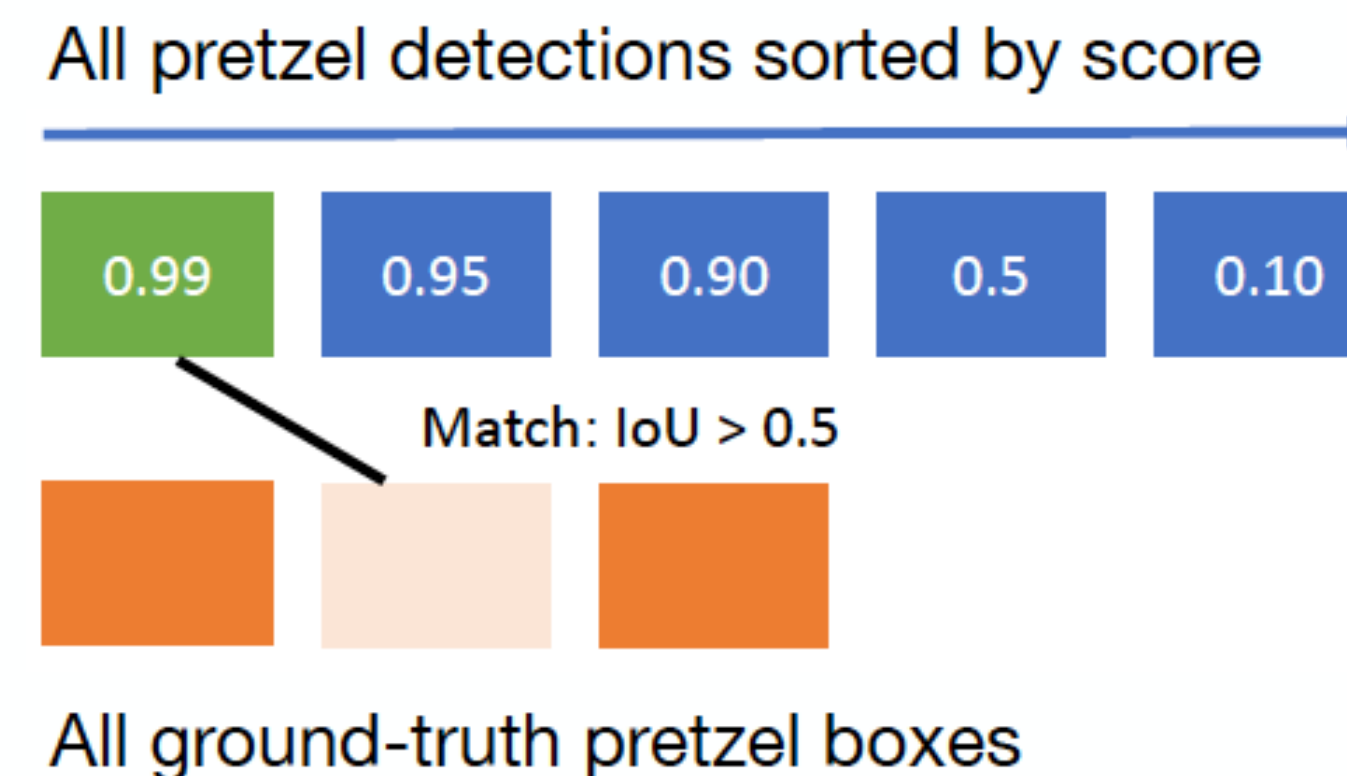
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative



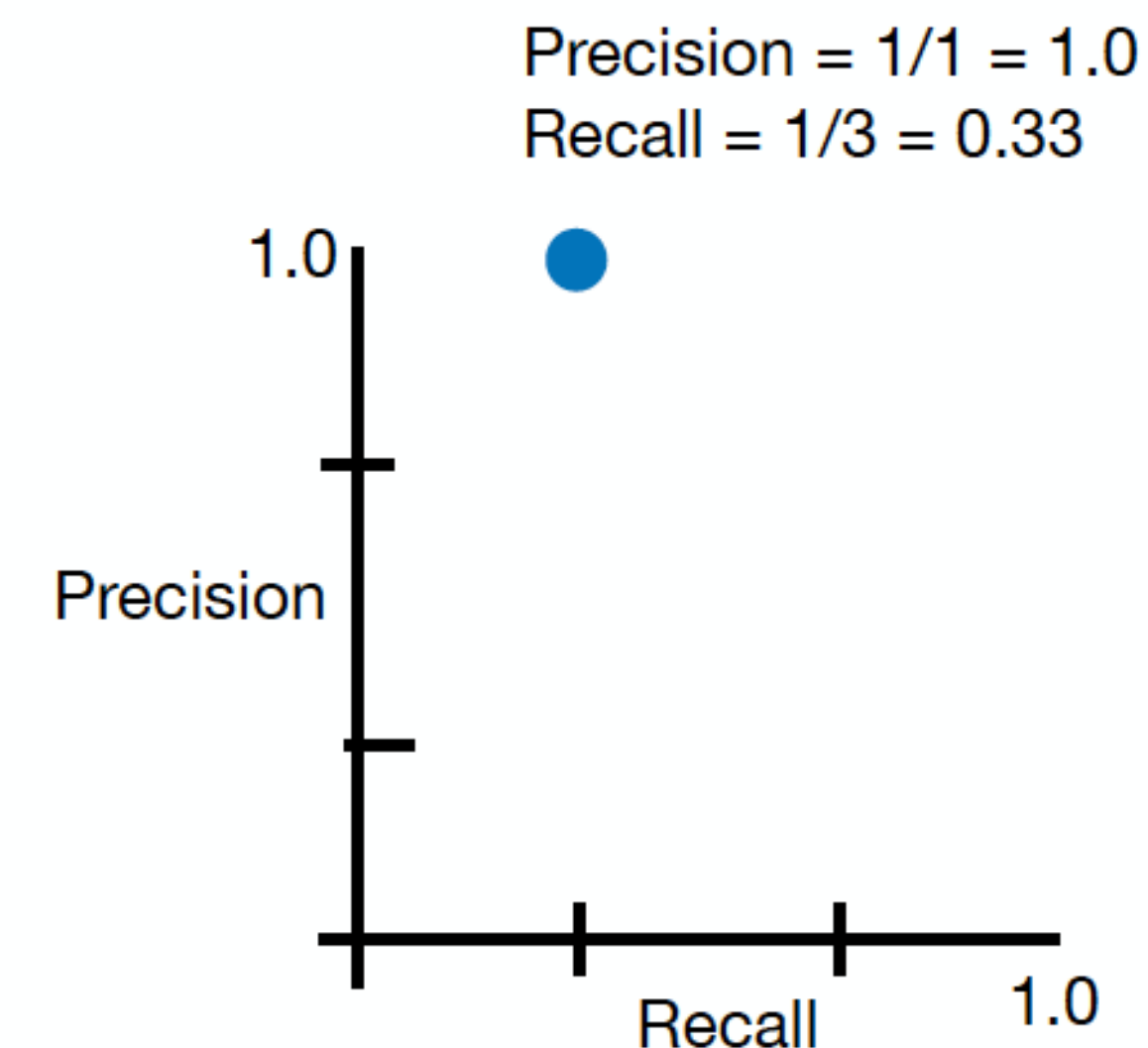
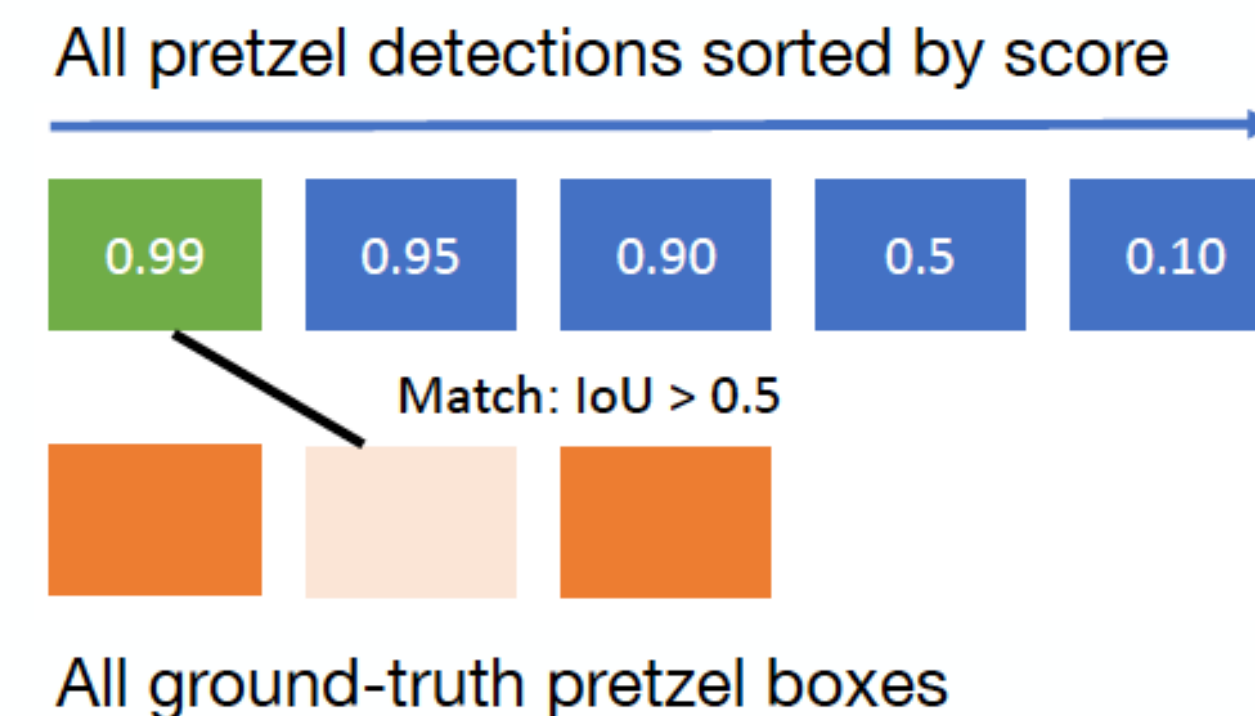


Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $IoU > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve

$$\text{Precision} = \frac{TP}{TP + FP}$$

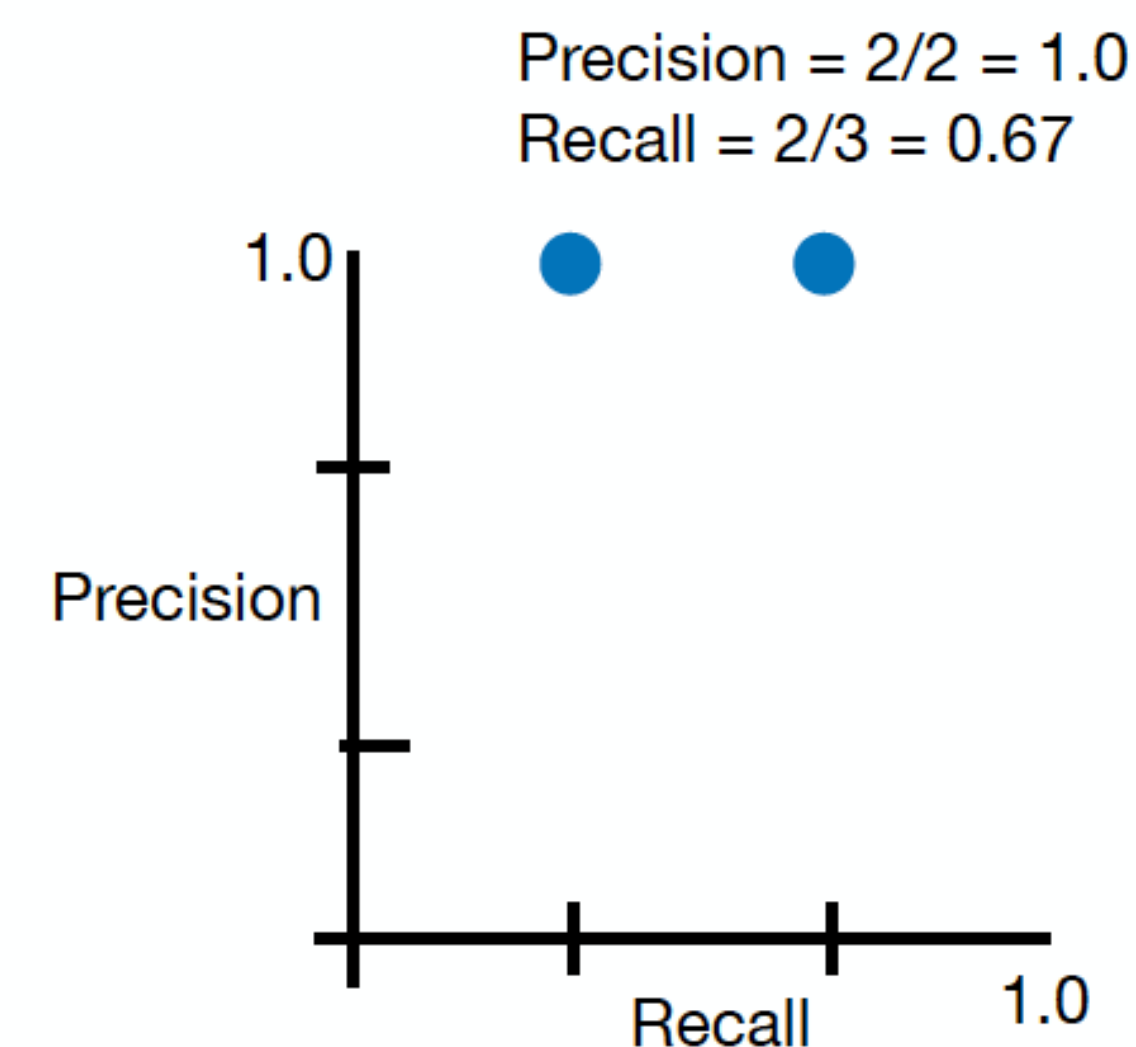
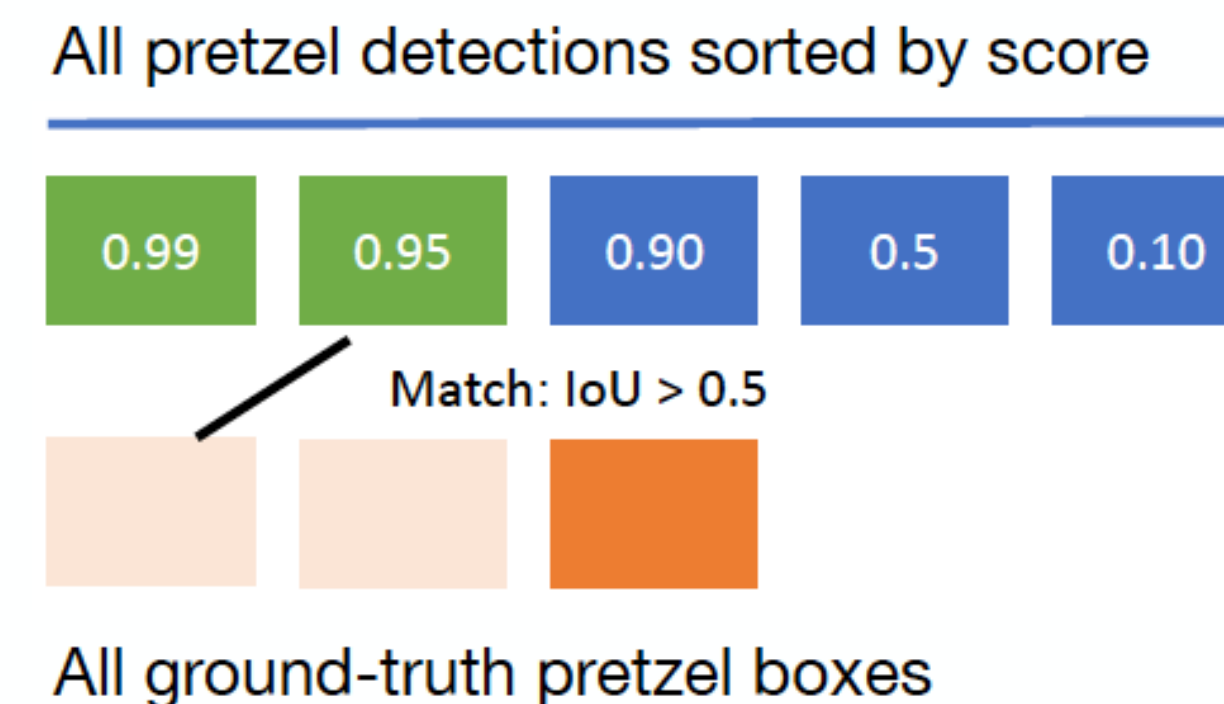
$$\text{Recall} = \frac{TP}{TP + FN}$$





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

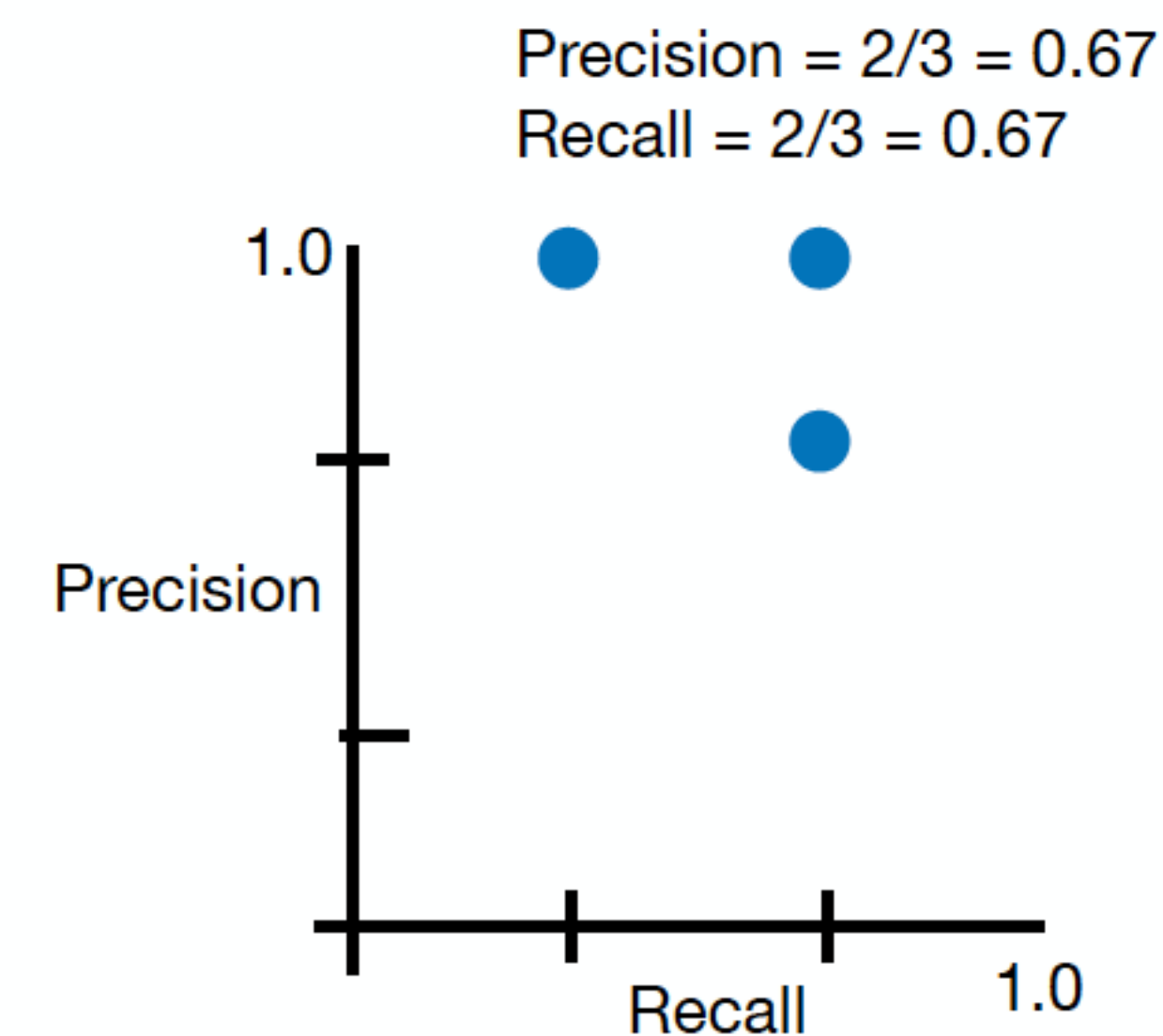
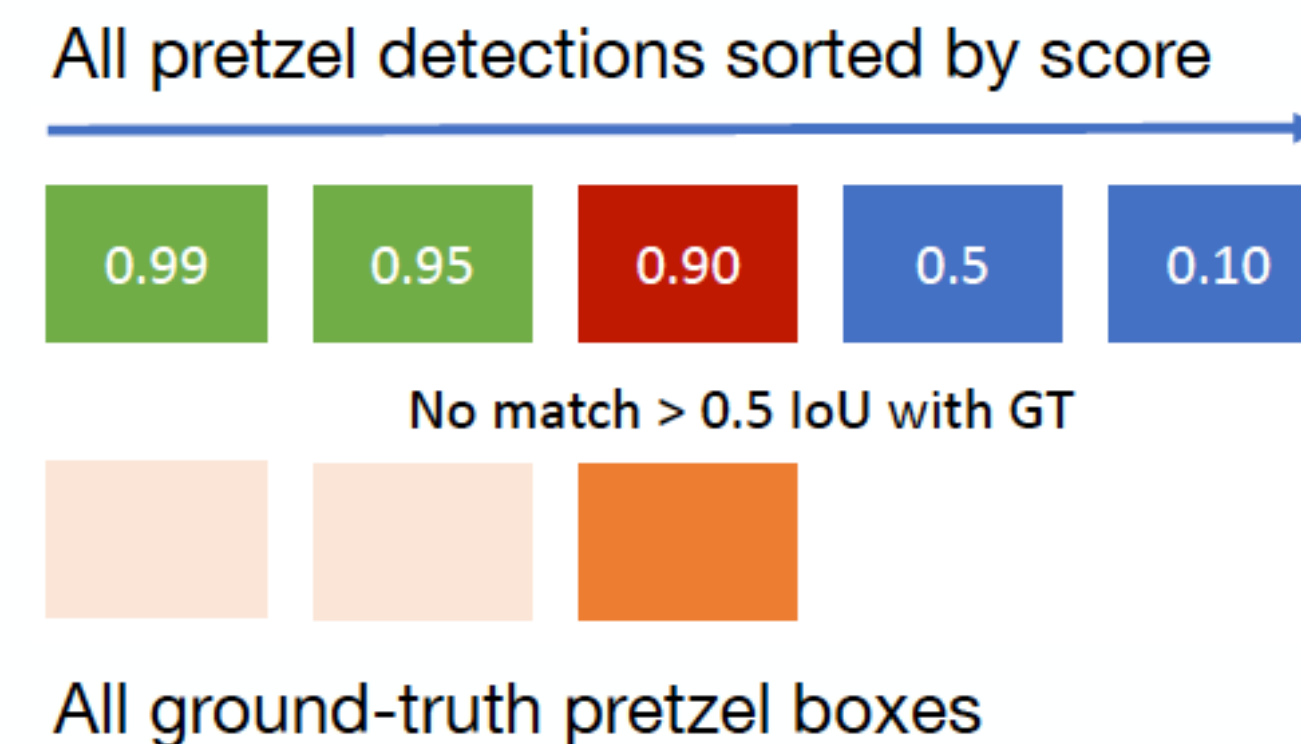
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

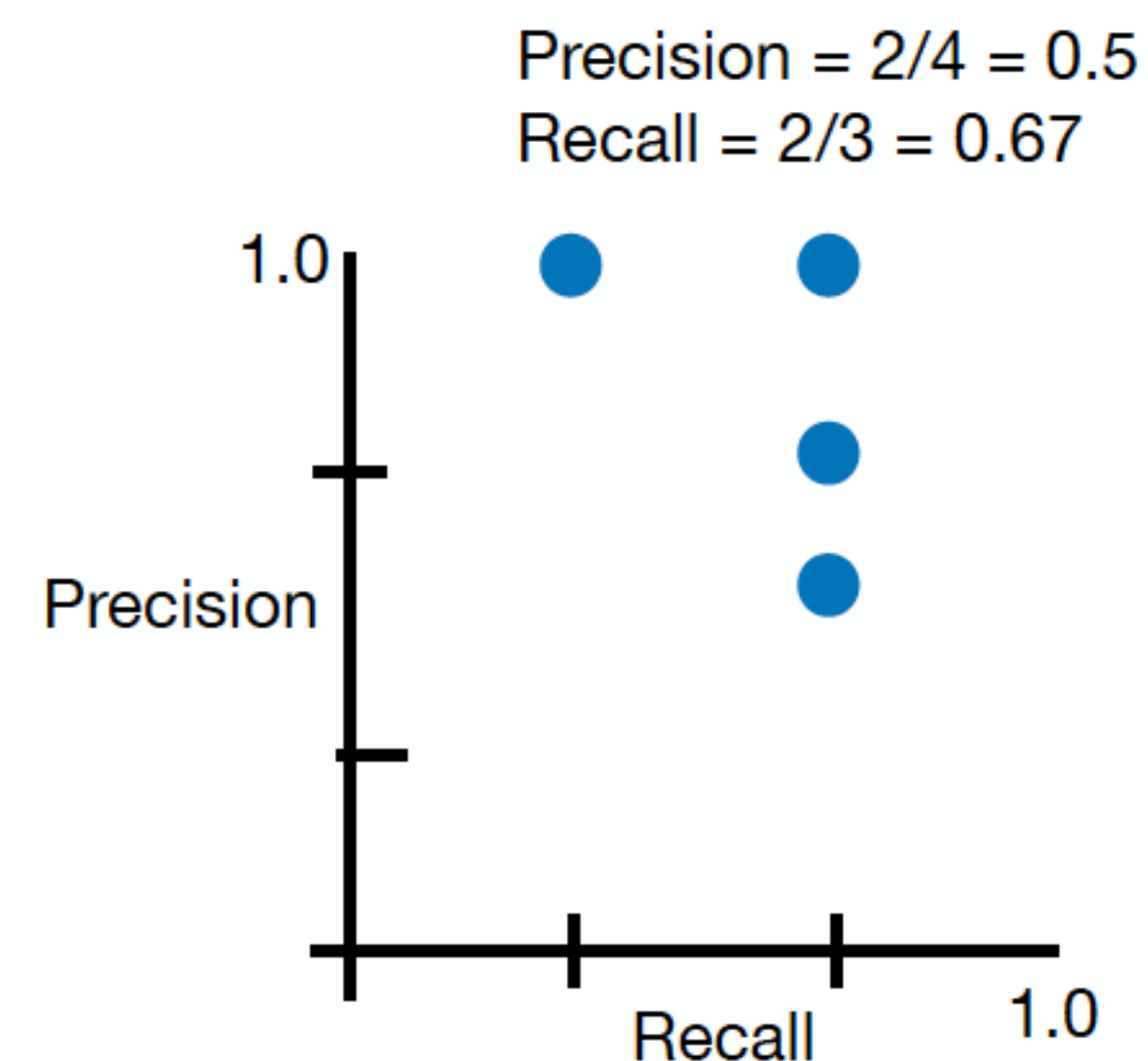
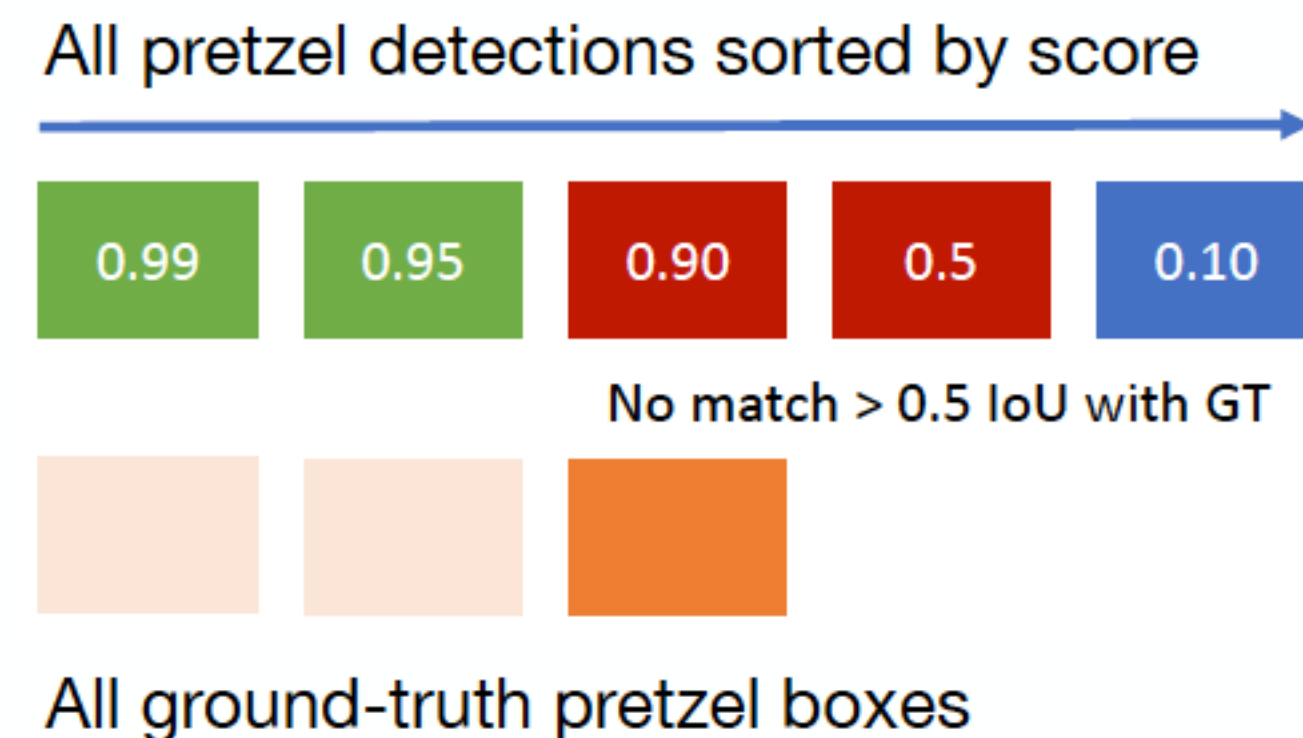
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

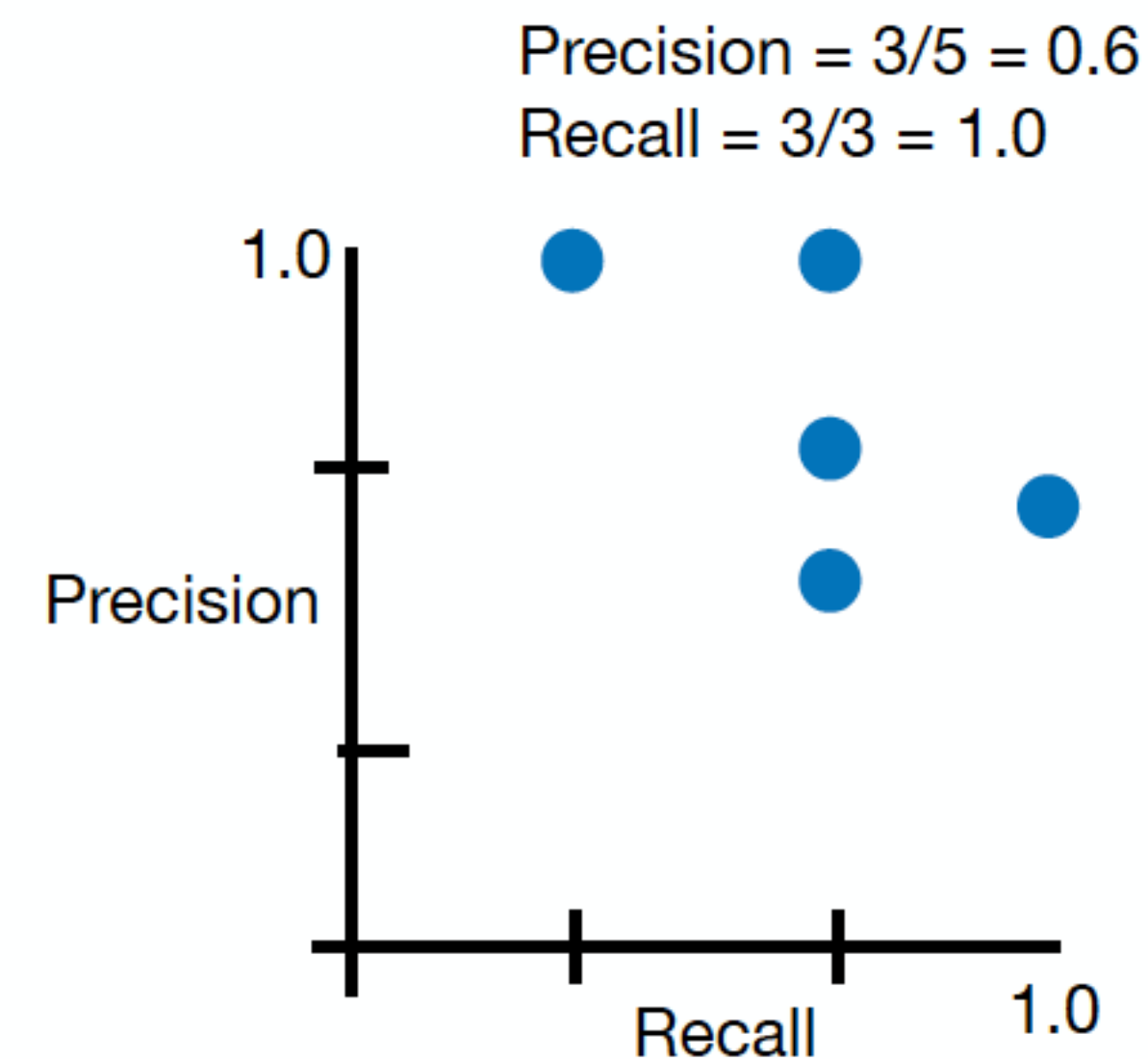
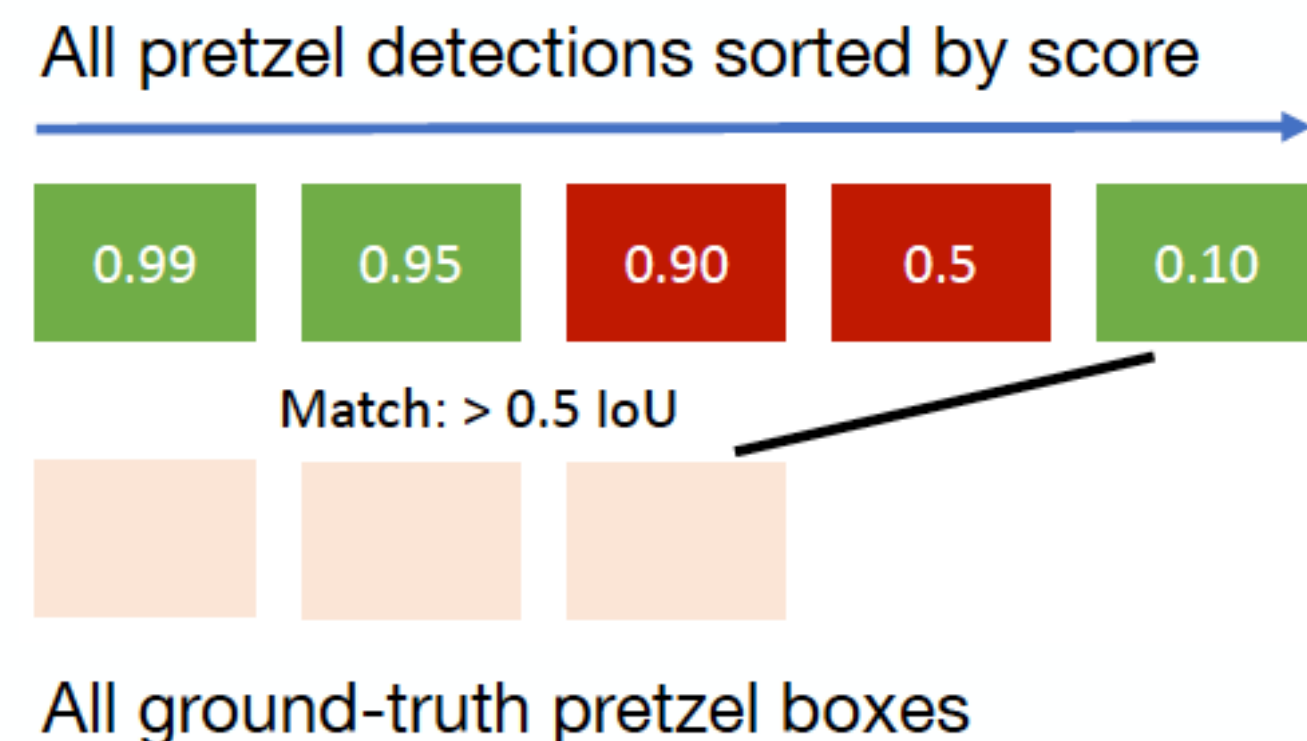
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

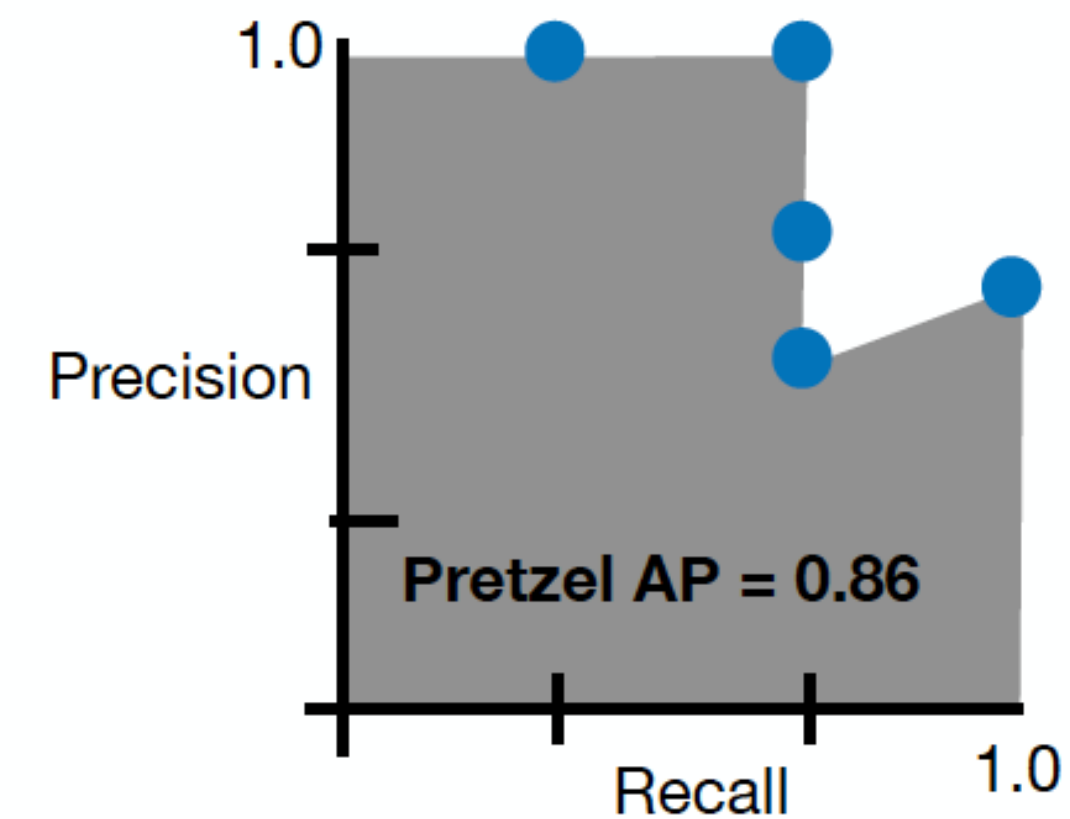
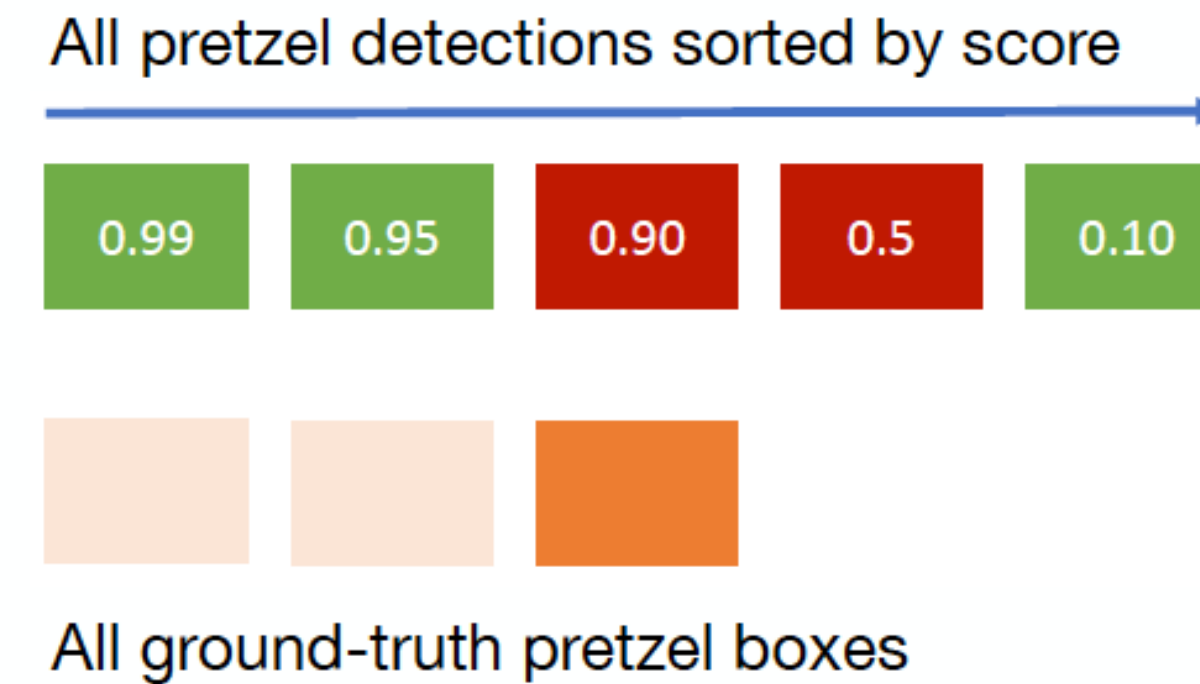
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve
 2. Average Precision (AP) = area under PR curve





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

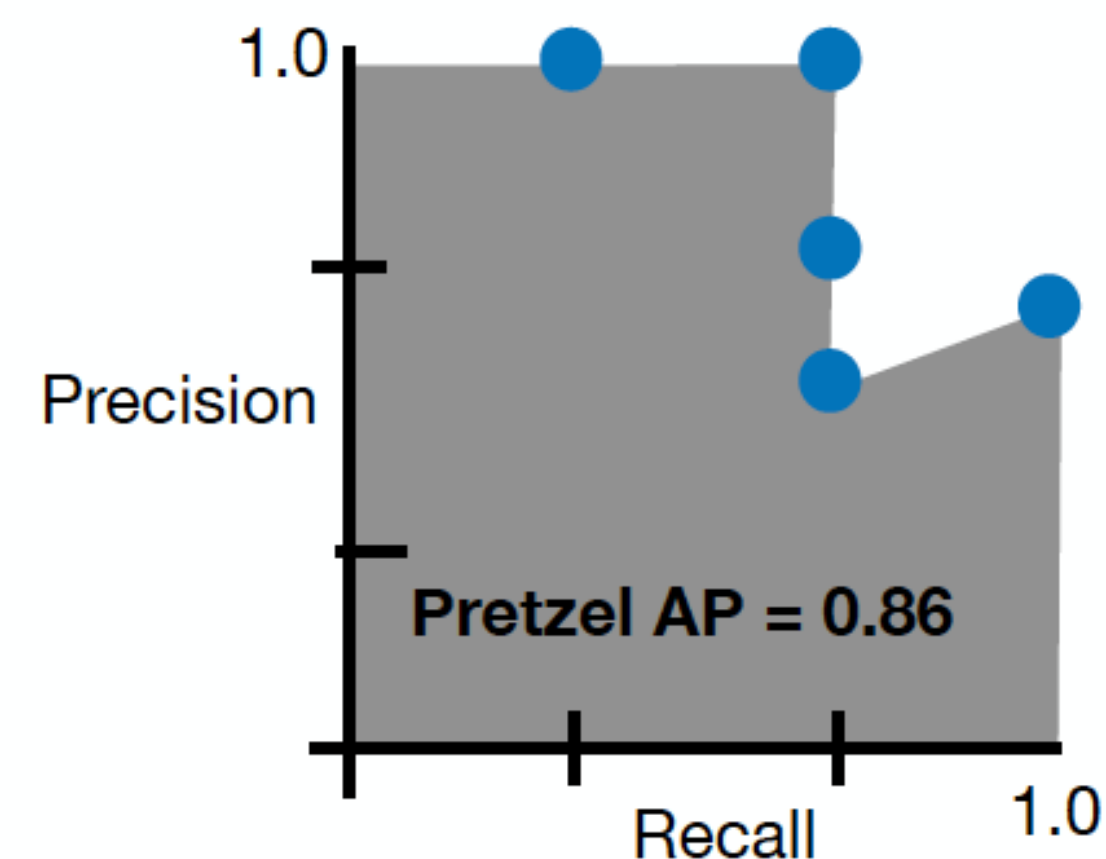
1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve
 2. Average Precision (AP) = area under PR curve

How to get AP = 1.0: Hit all GT boxes with IoU > 0.5, and have no “false positive” detections ranked above any “true positives”

All pretzel detections sorted by score



All ground-truth pretzel boxes





Evaluating Object Detectors: Mean Average Precision (mAP) and P-R Curve

1. Run object detector on all test images (with NMS)
2. For each category, compute Average Precision (AP) = area under Precision vs Recall Curve
 1. For each detection (highest score to lowest score)
 1. If it matches some GT box with $\text{IoU} > 0.5$, mark it as positive and eliminate the GT
 2. Otherwise mark it as negative
 3. Plot a point on PR curve
 2. Average Precision (AP) = area under PR curve
3. Mean Average Precision (mAP) = average of AP for each category

Flipz AP = 0.60
Hershey's AP = 0.85
Reese's AP = 0.81
mAP@0.5 = 0.75

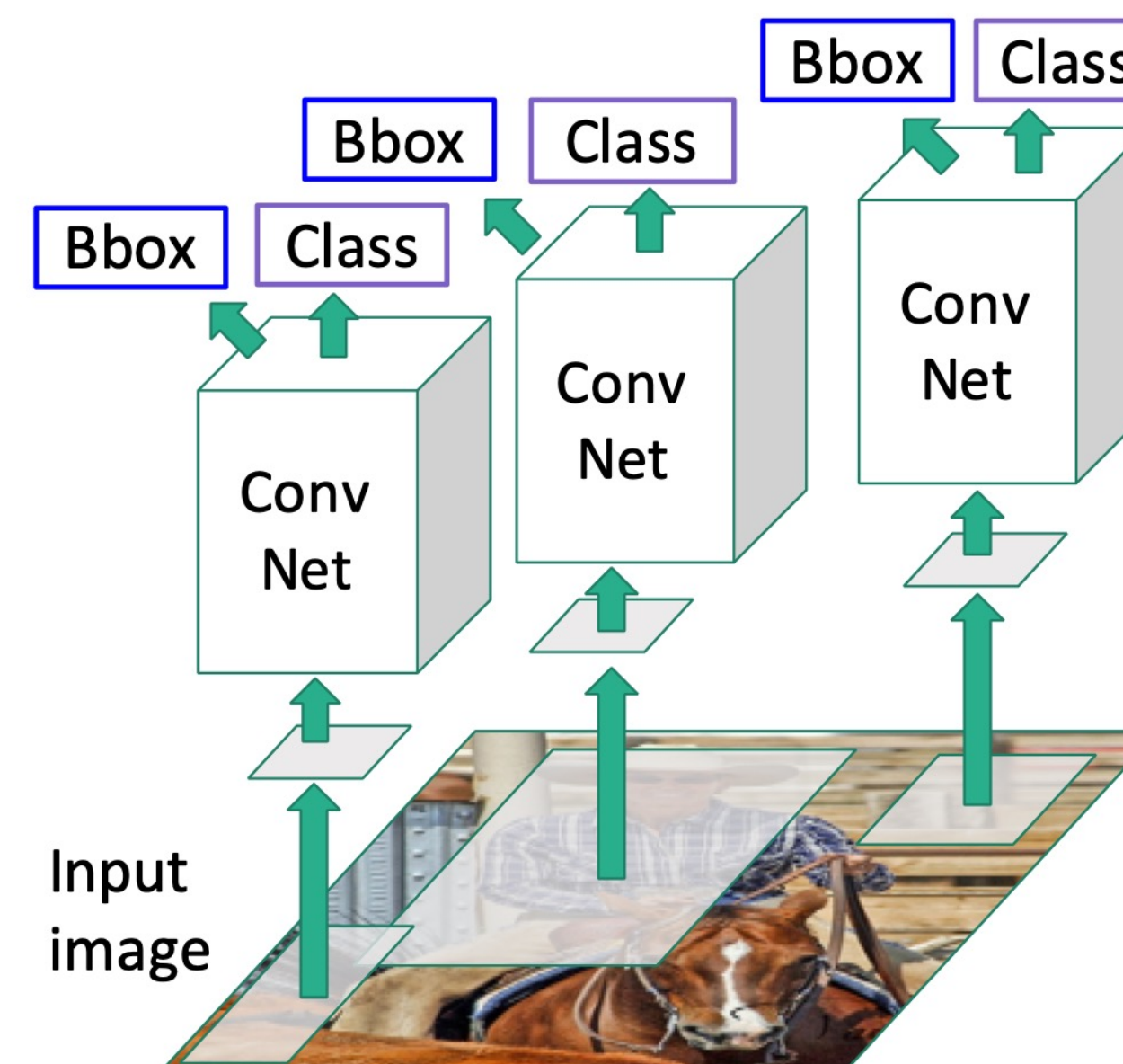
“COCO Evaluator”



Fast R-CNN

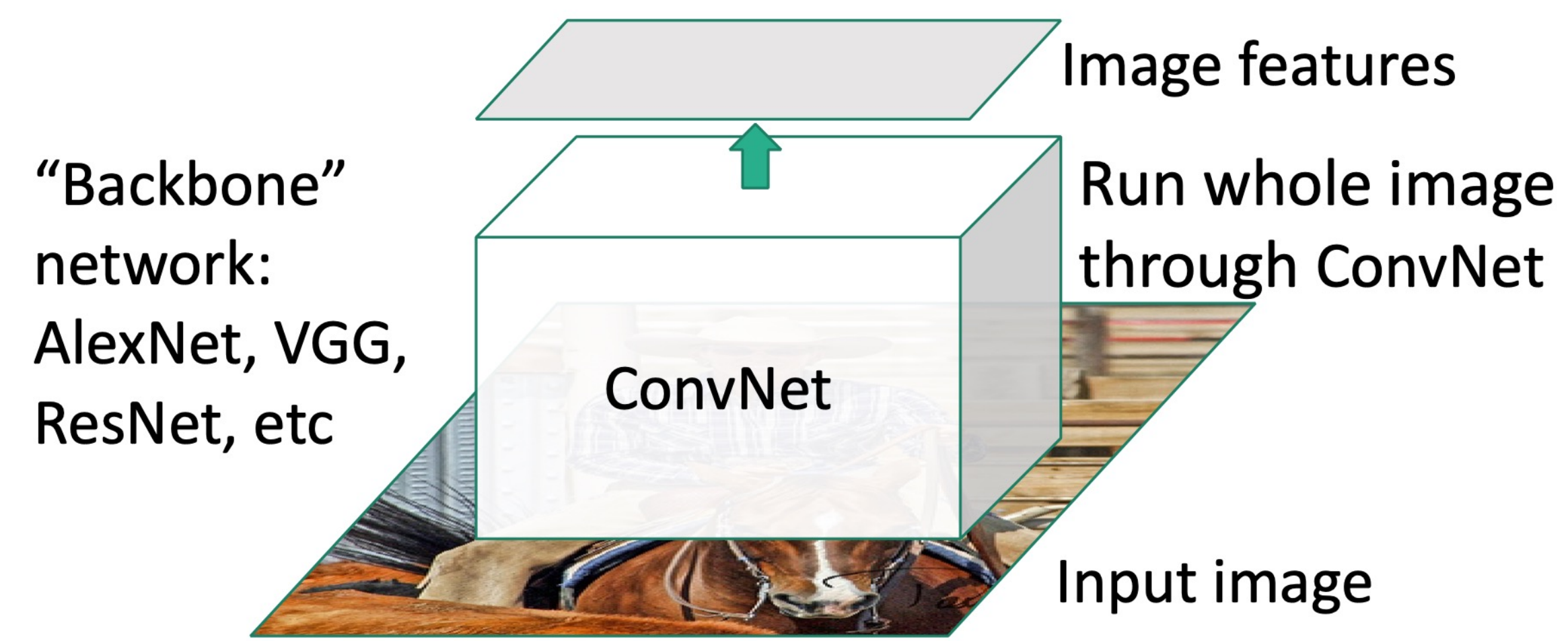


“Slow” R-CNN
Process each region
independently

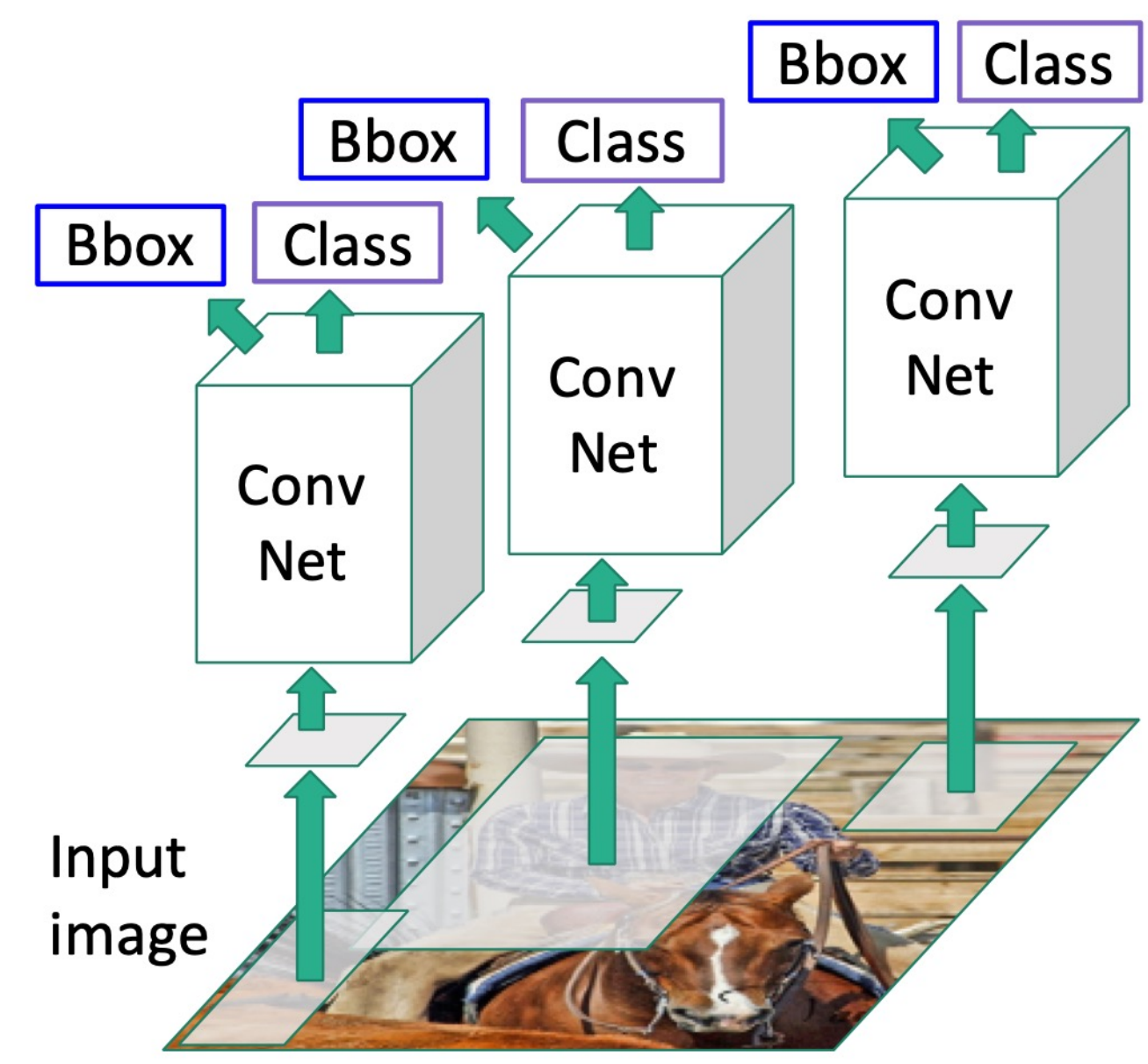




Fast R-CNN



“Slow” R-CNN
Process each region independently

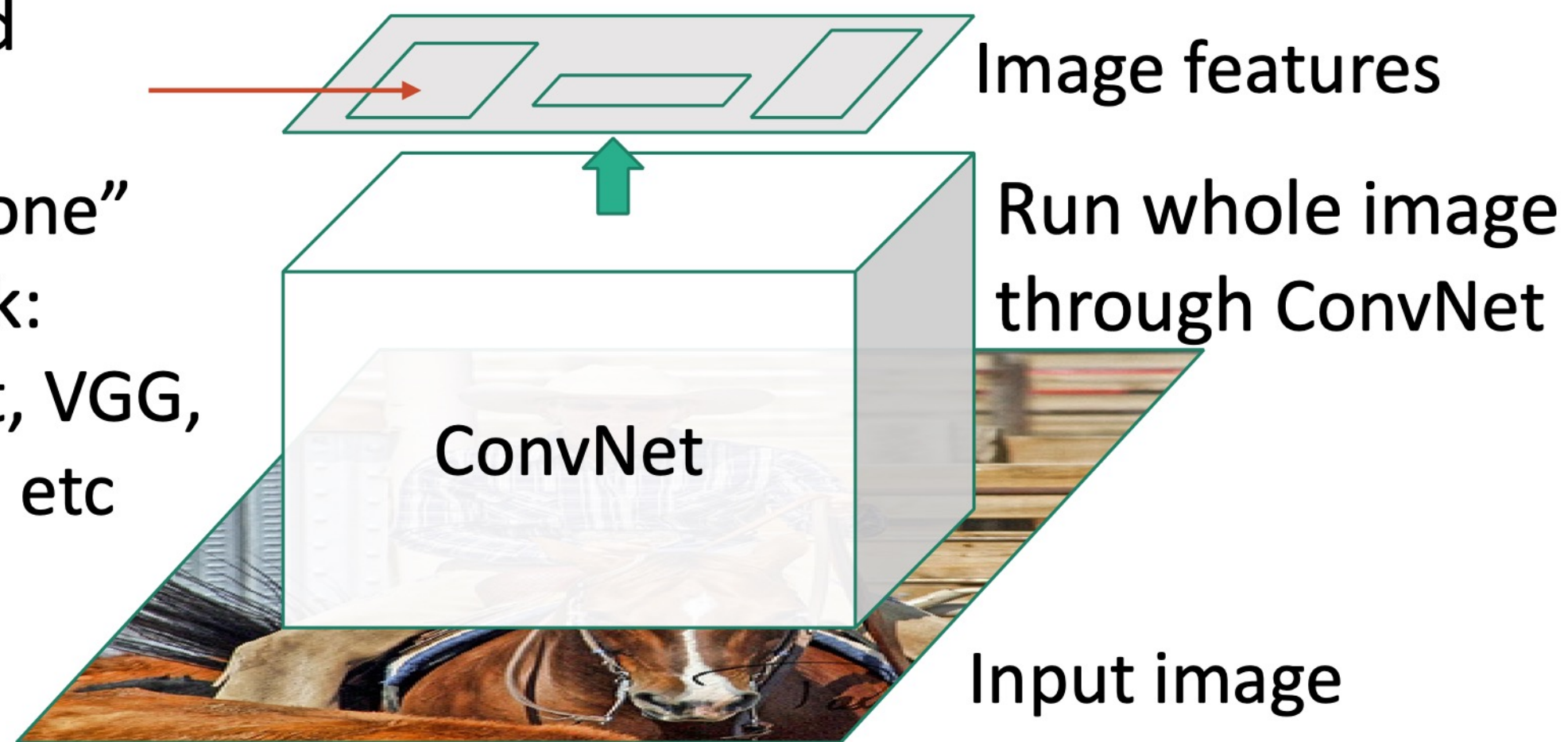




Fast R-CNN

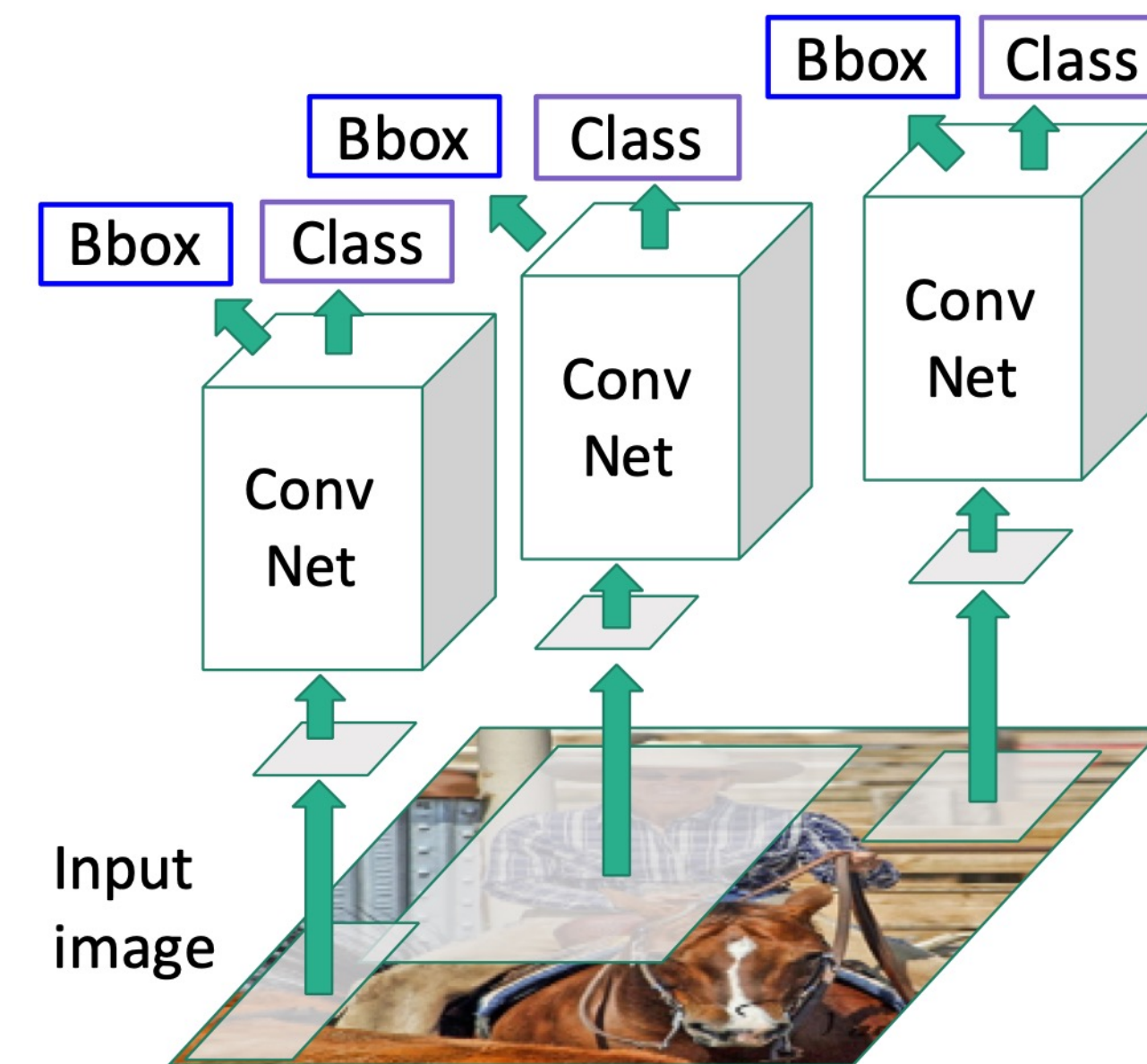
Regions of Interest (RoIs) from a proposal method

“Backbone” network:
AlexNet, VGG,
ResNet, etc



“Slow” R-CNN

Process each region independently

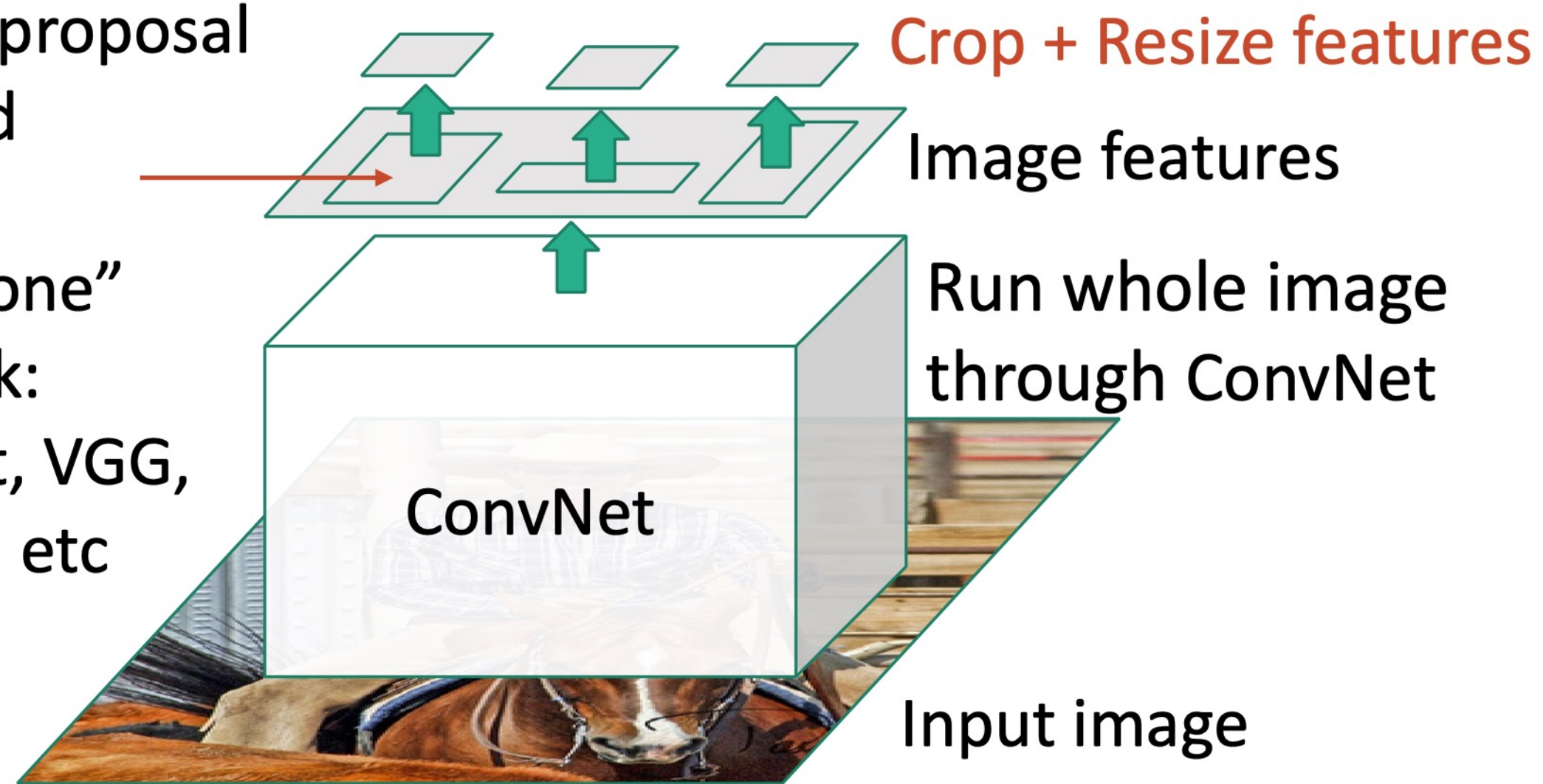




Fast R-CNN

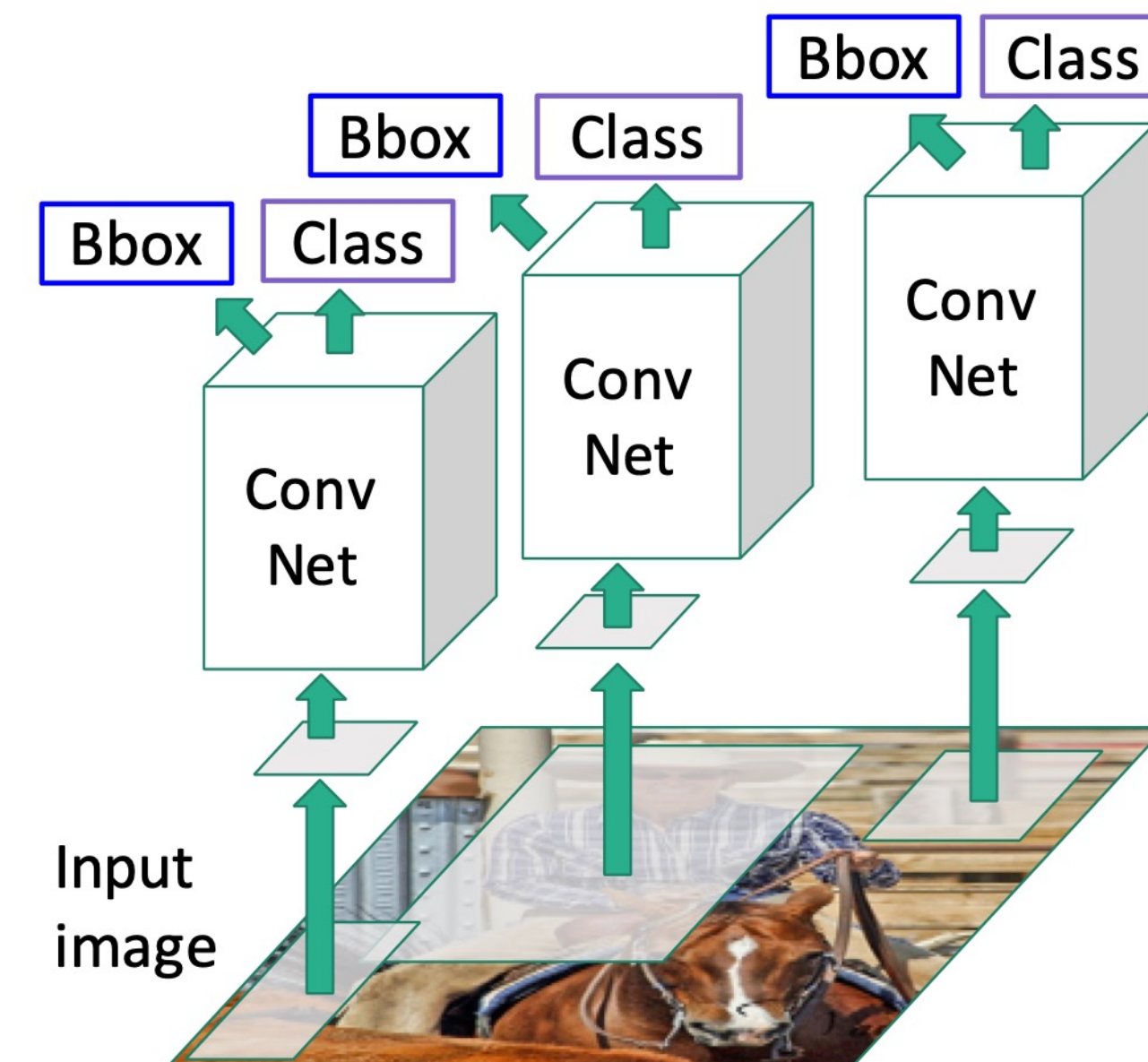
Regions of Interest (Rois) from a proposal method

“Backbone” network:
AlexNet, VGG, ResNet, etc



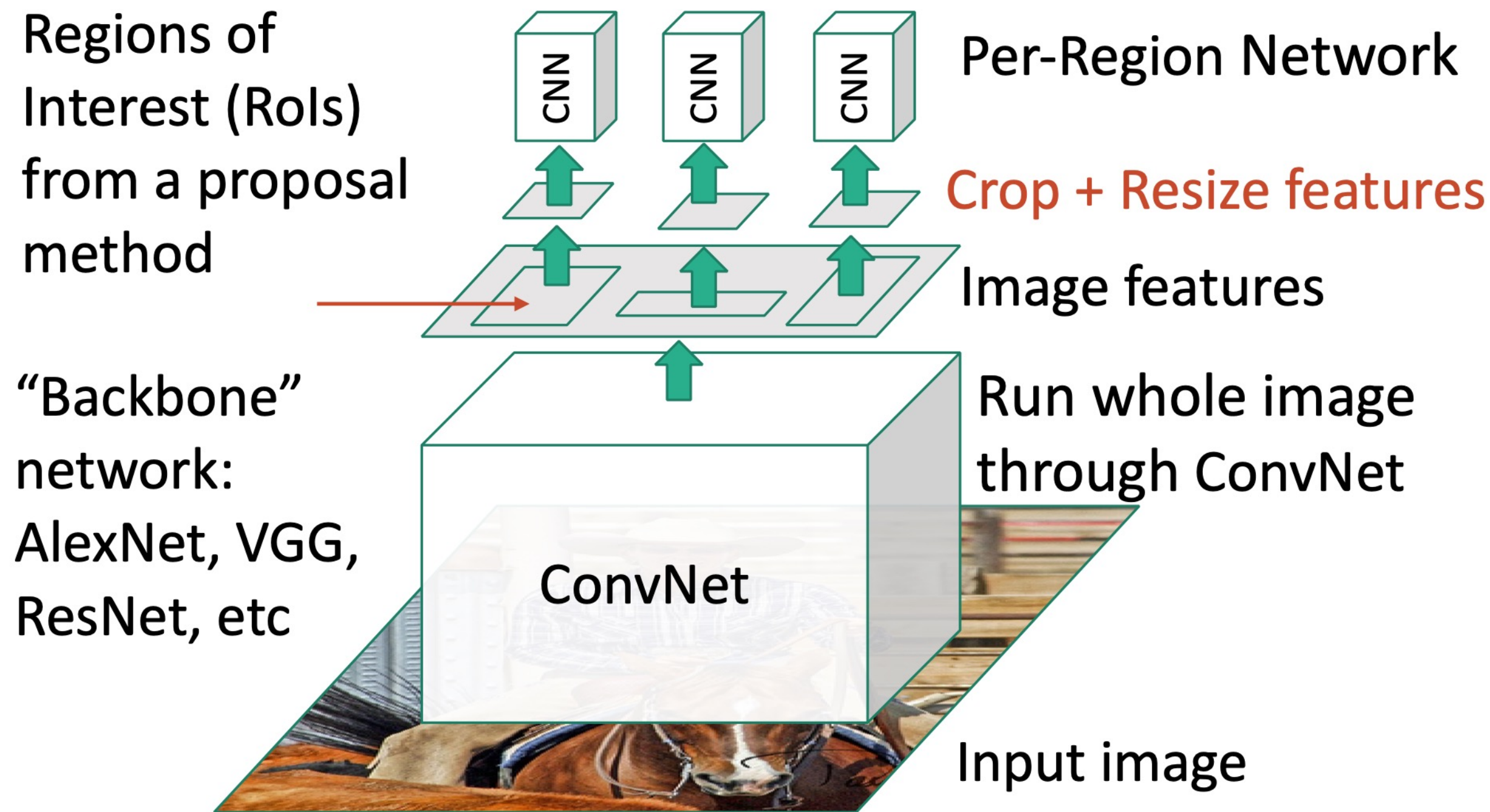
“Slow” R-CNN

Process each region independently

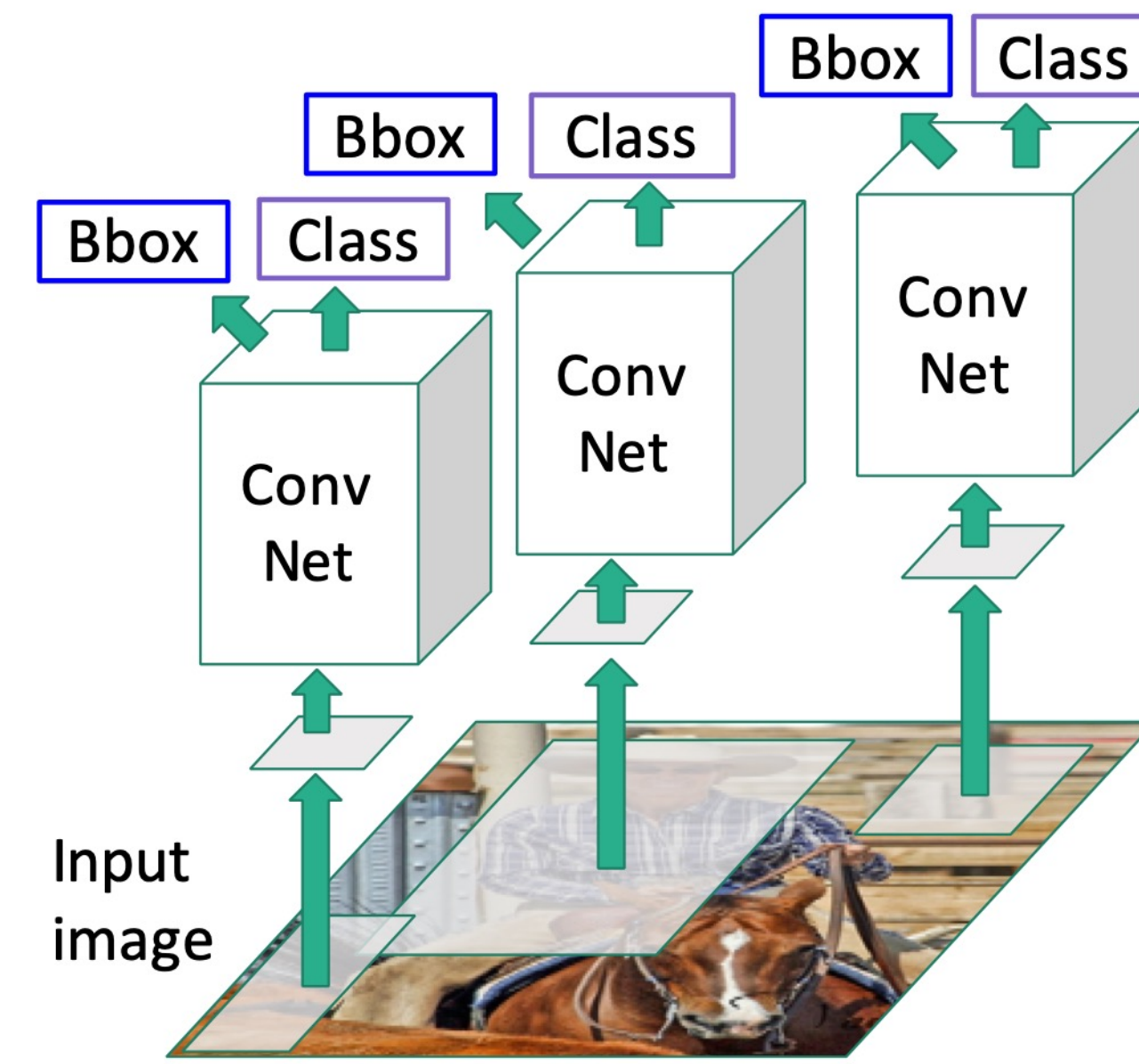




Fast R-CNN

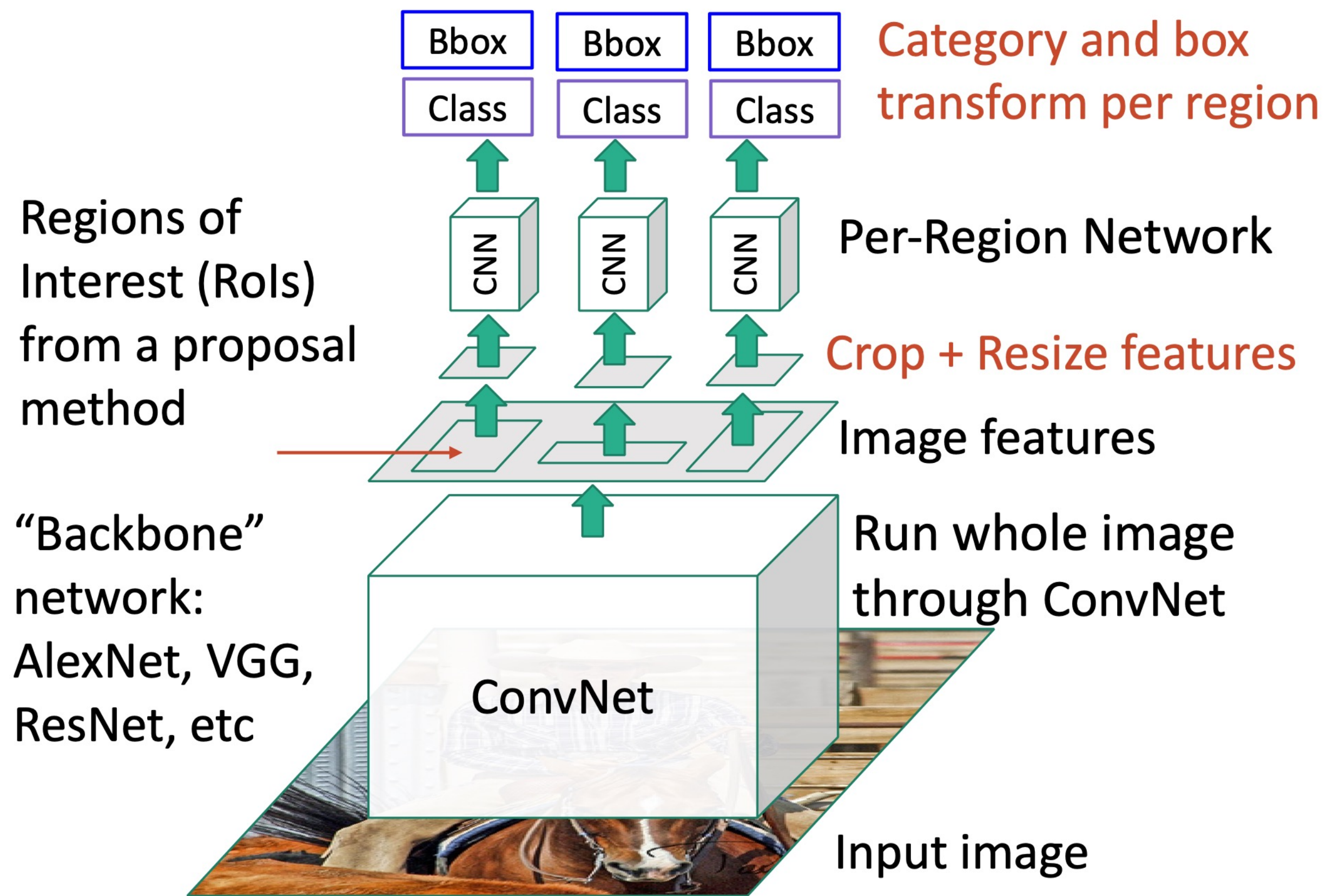


“Slow” R-CNN
Process each region independently

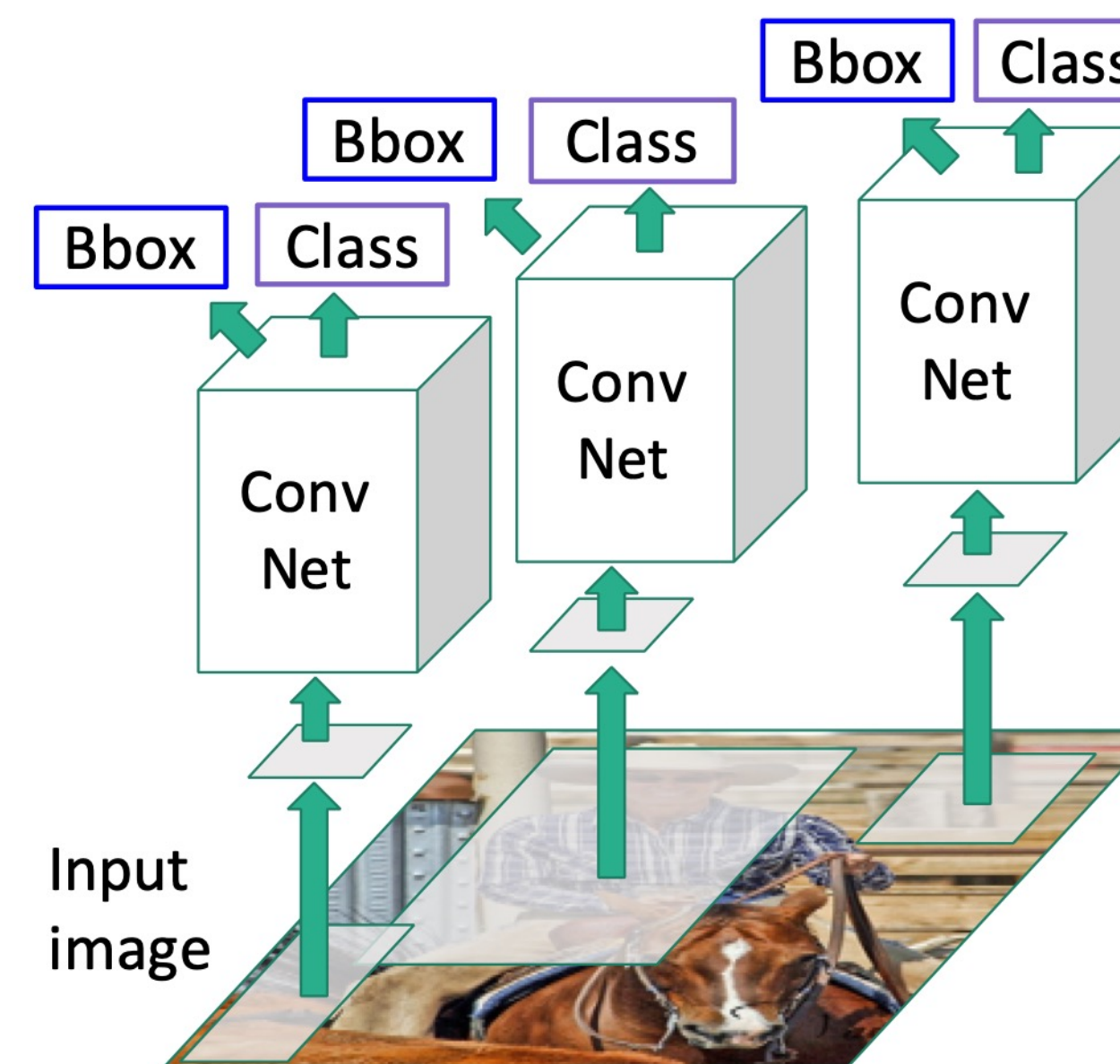




Fast R-CNN

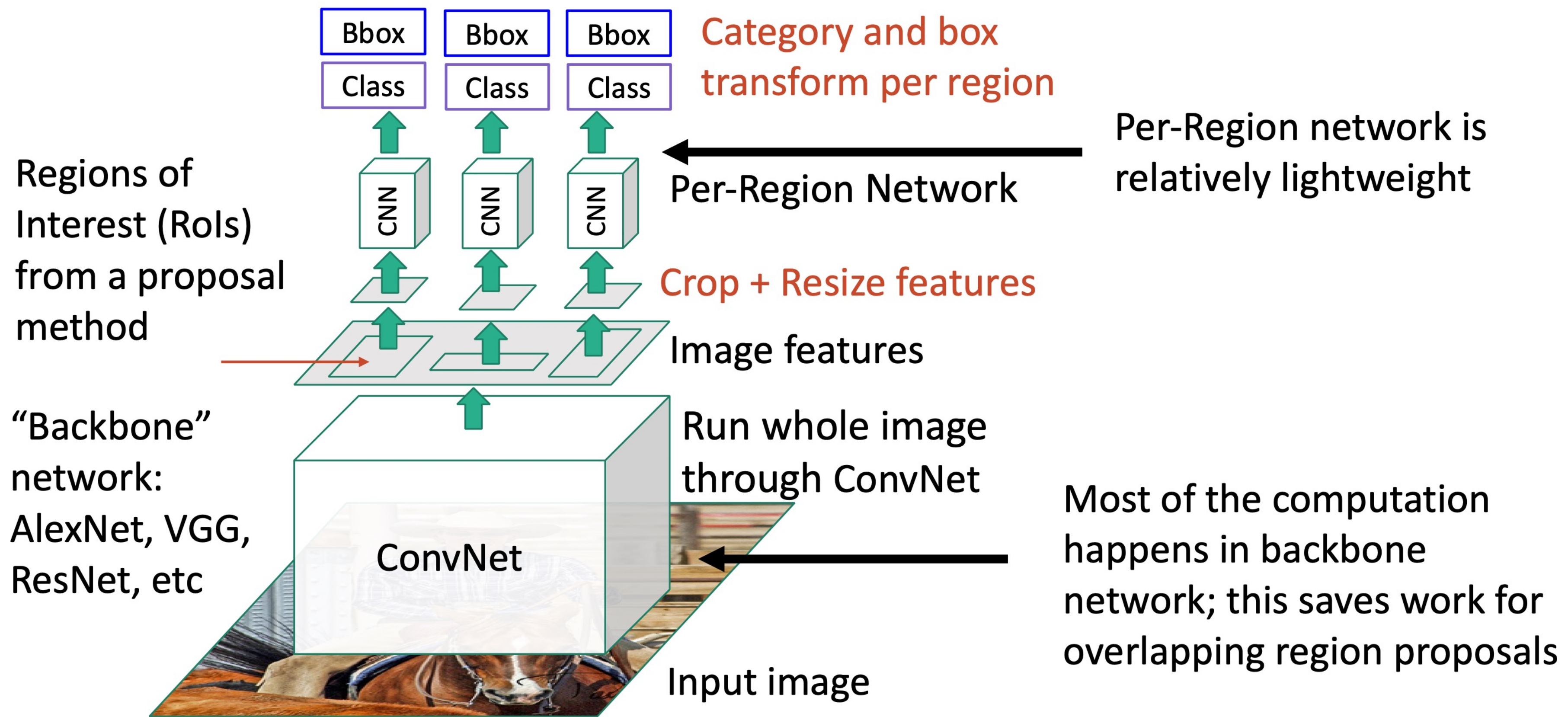


“Slow” R-CNN
Process each region independently



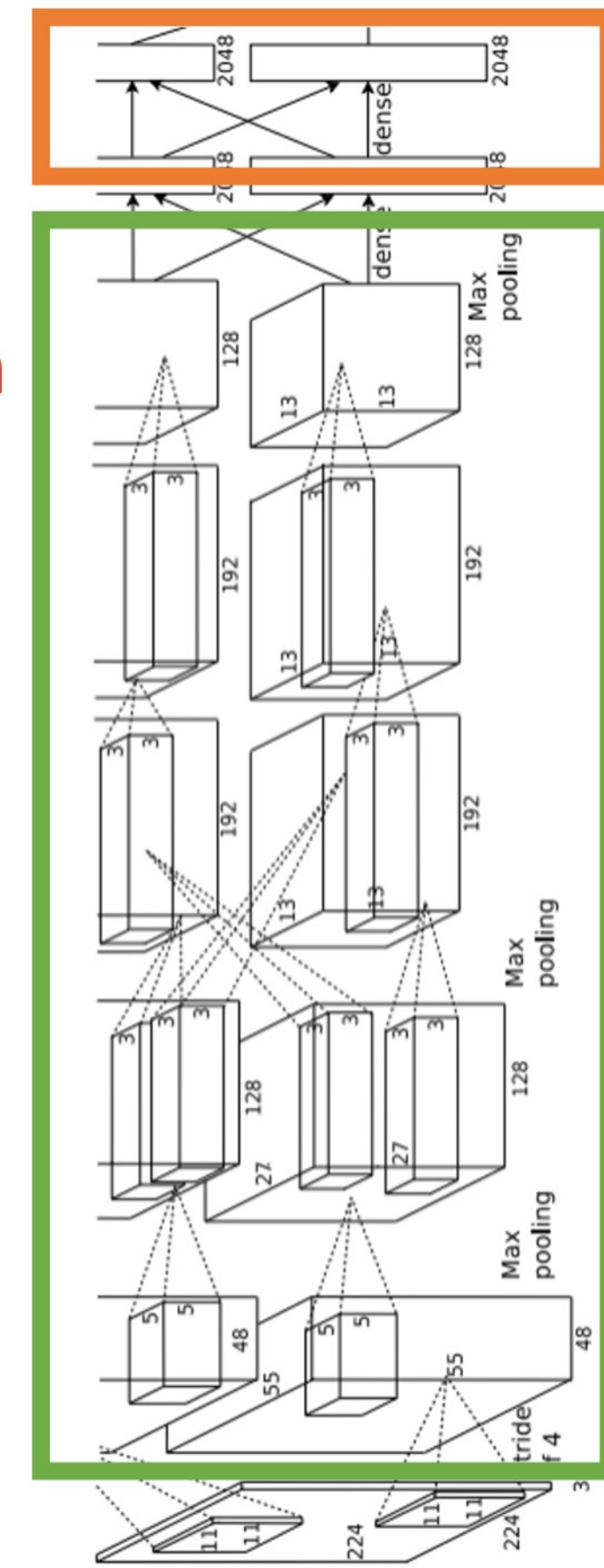
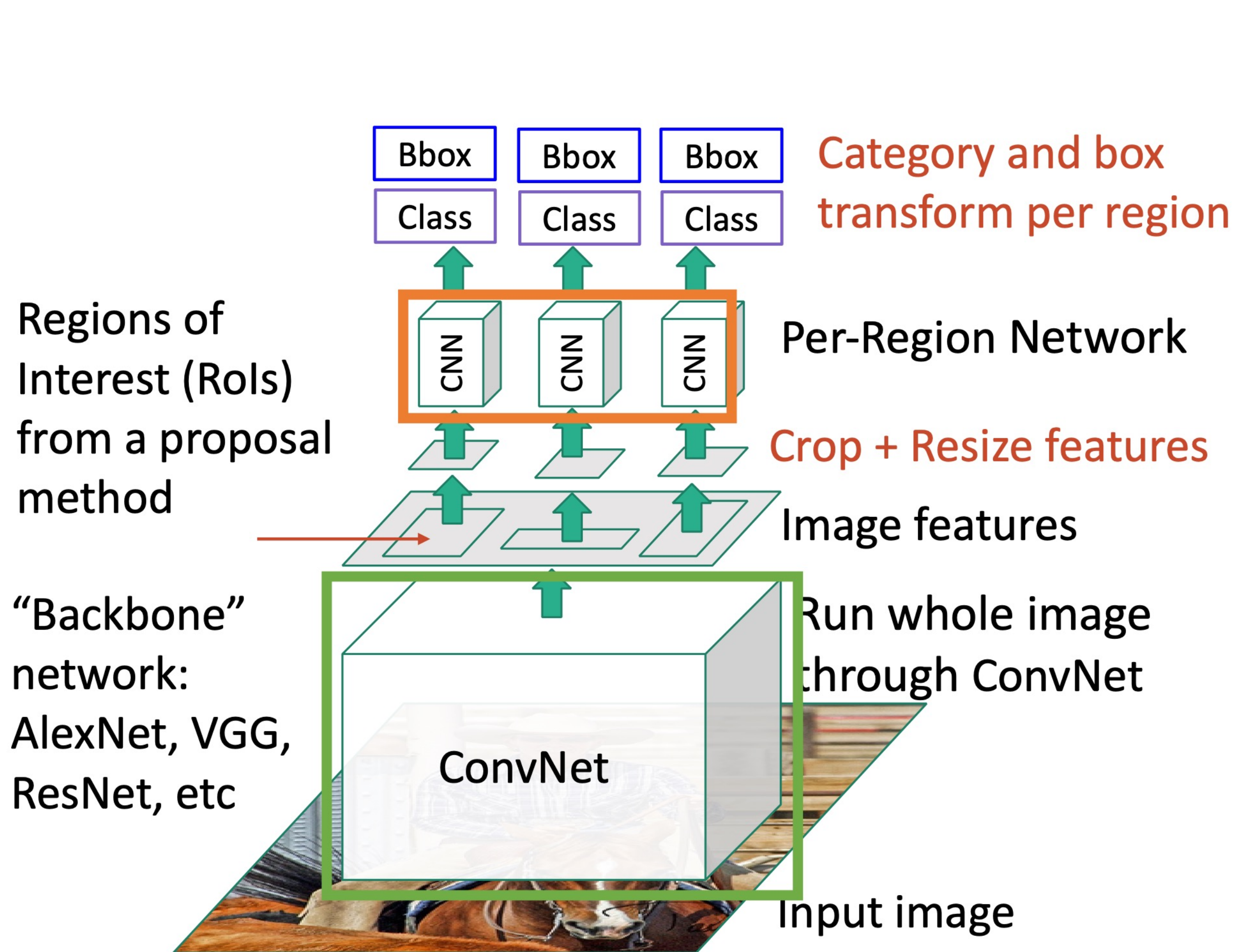


Fast R-CNN





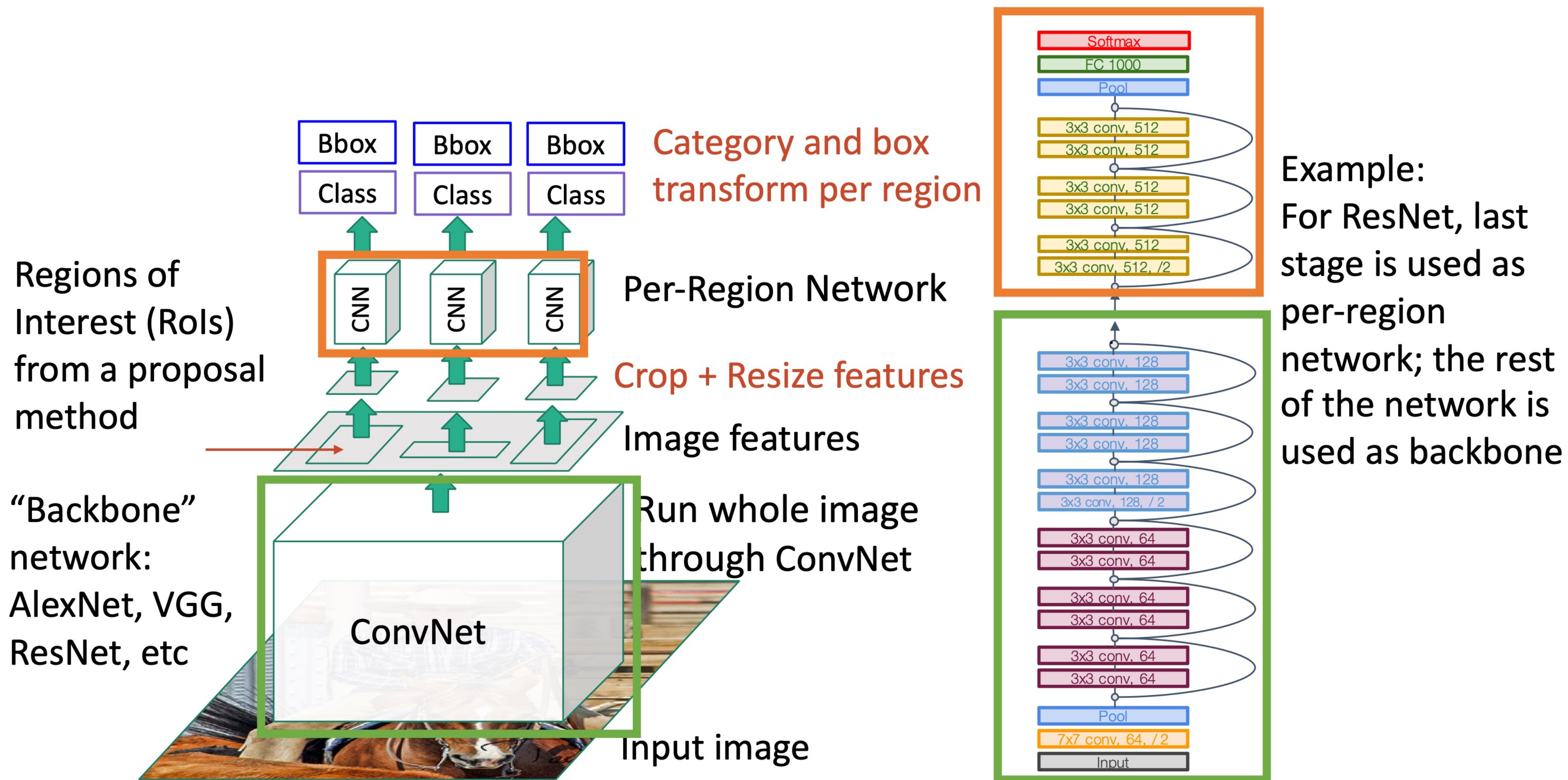
Fast R-CNN



Example:
When using AlexNet for detection, five conv layers are used for backbone and two FC layers are used for per-region network

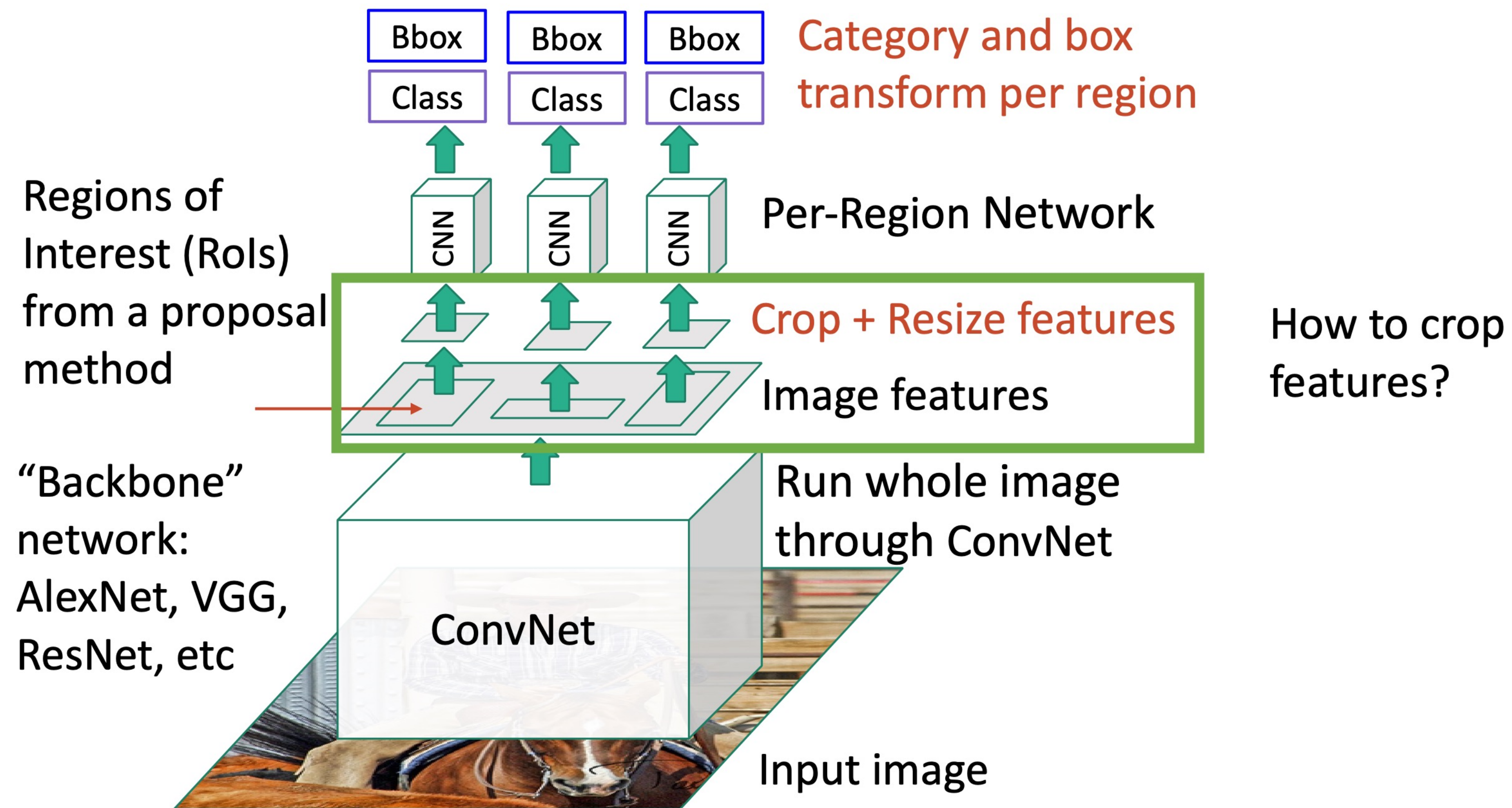


Fast R-CNN



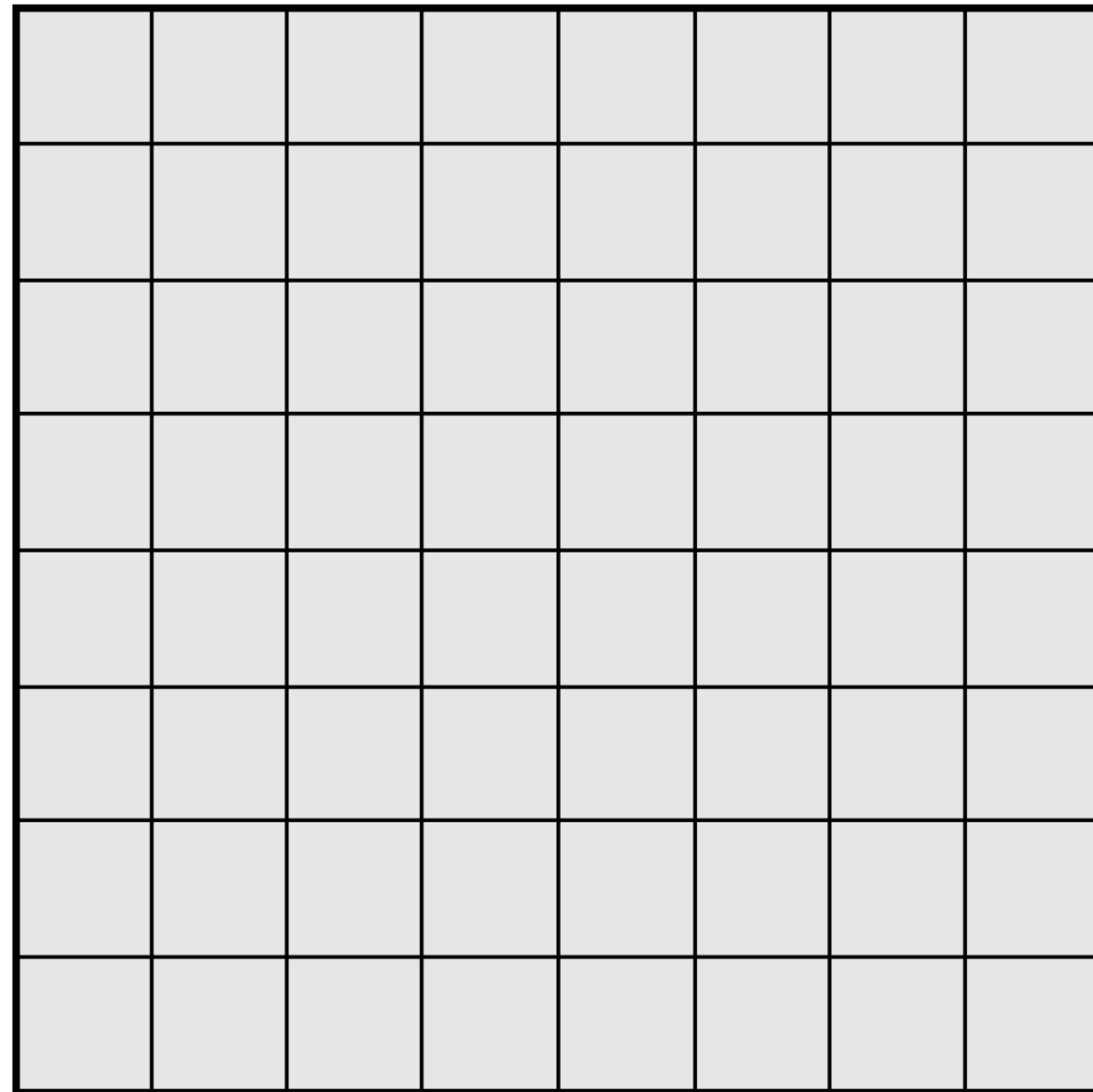


Fast R-CNN





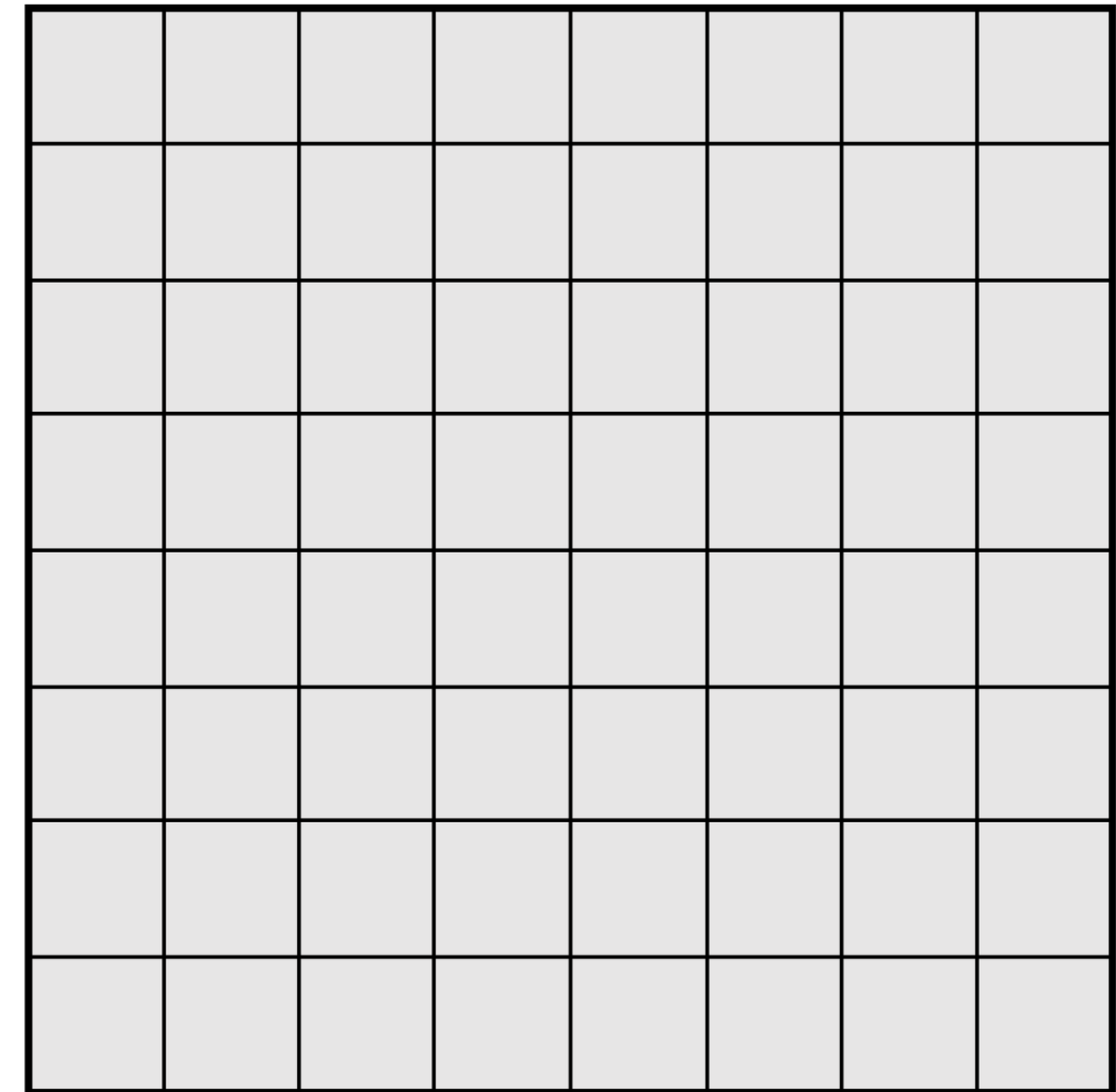
Recall: Receptive Fields



Input Image: 8 x 8

Every position in the output feature map depends on a 3x3 receptive field in the input

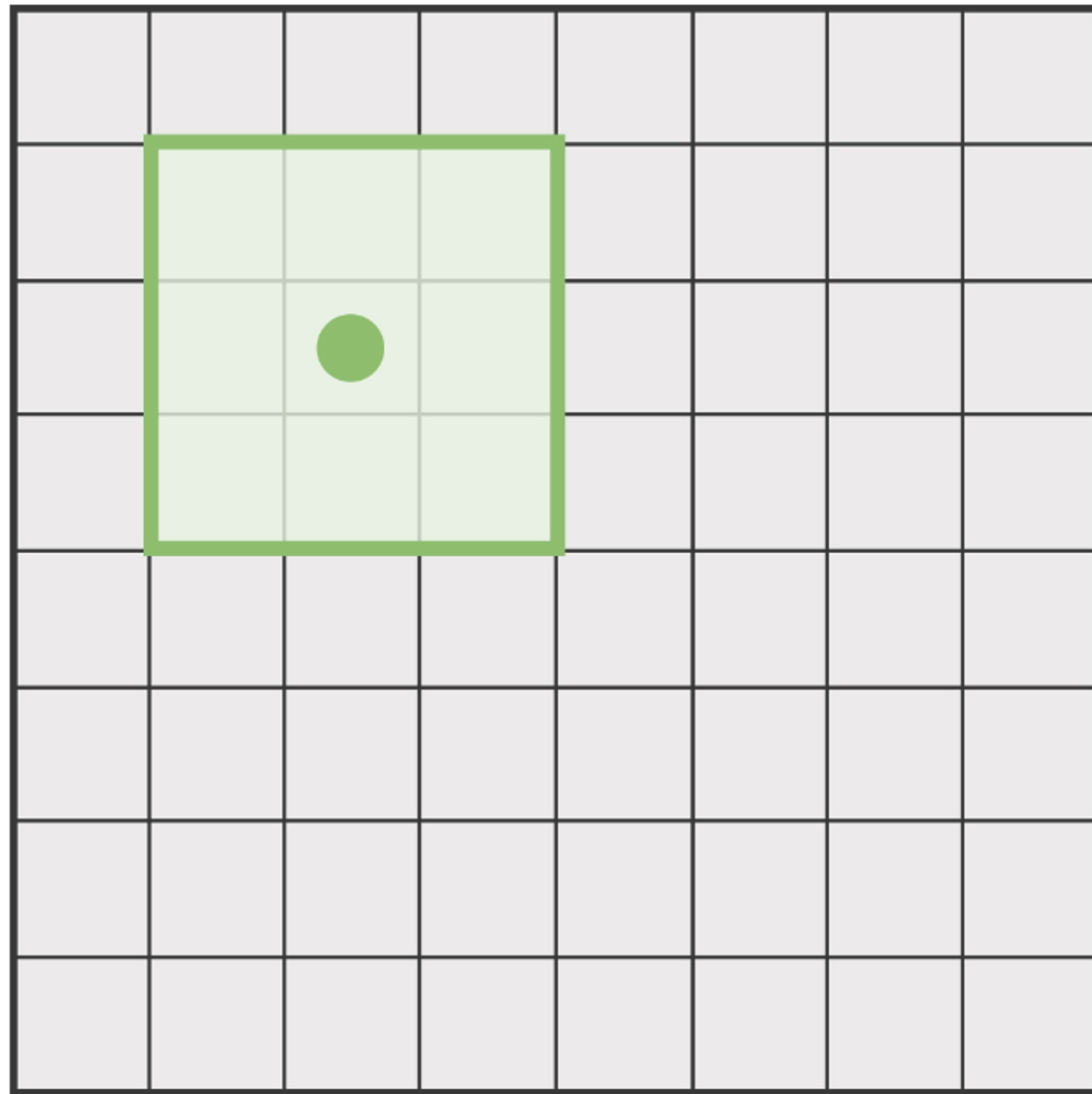
3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8



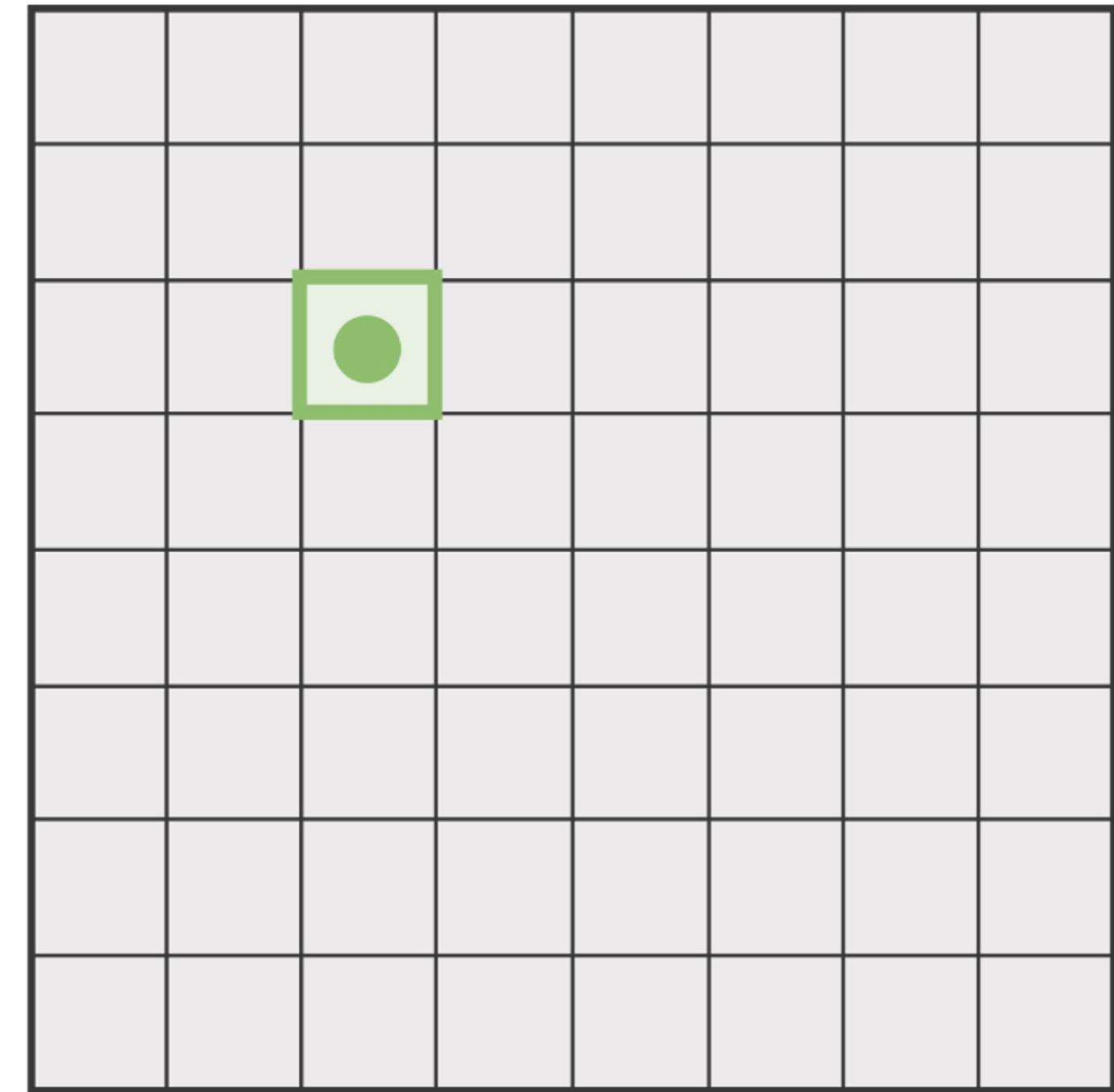
Recall: Receptive Fields



Input Image: 8 x 8

Every position in the output feature map depends on a 3x3 receptive field in the input

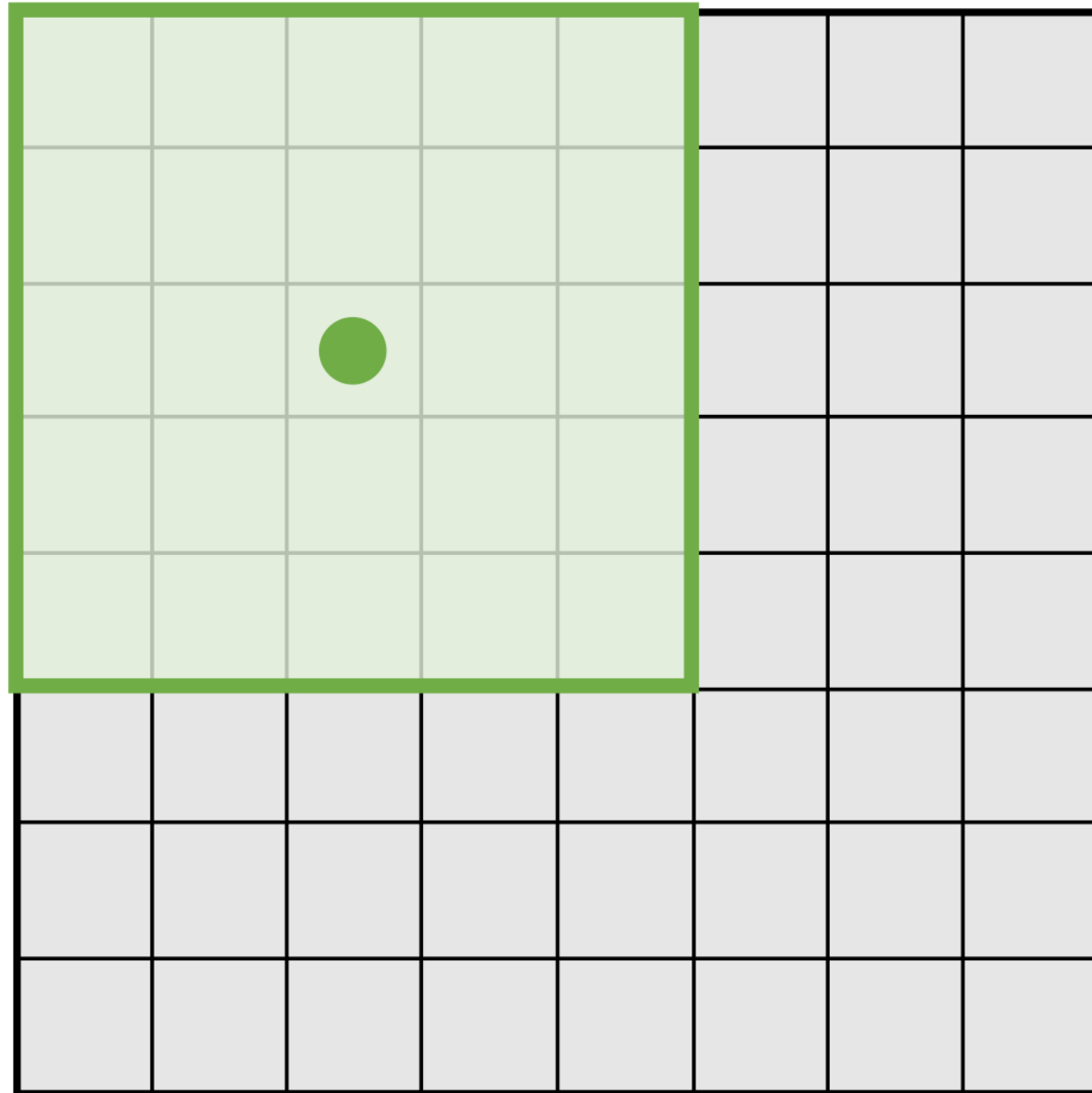
3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8



Recall: Receptive Fields

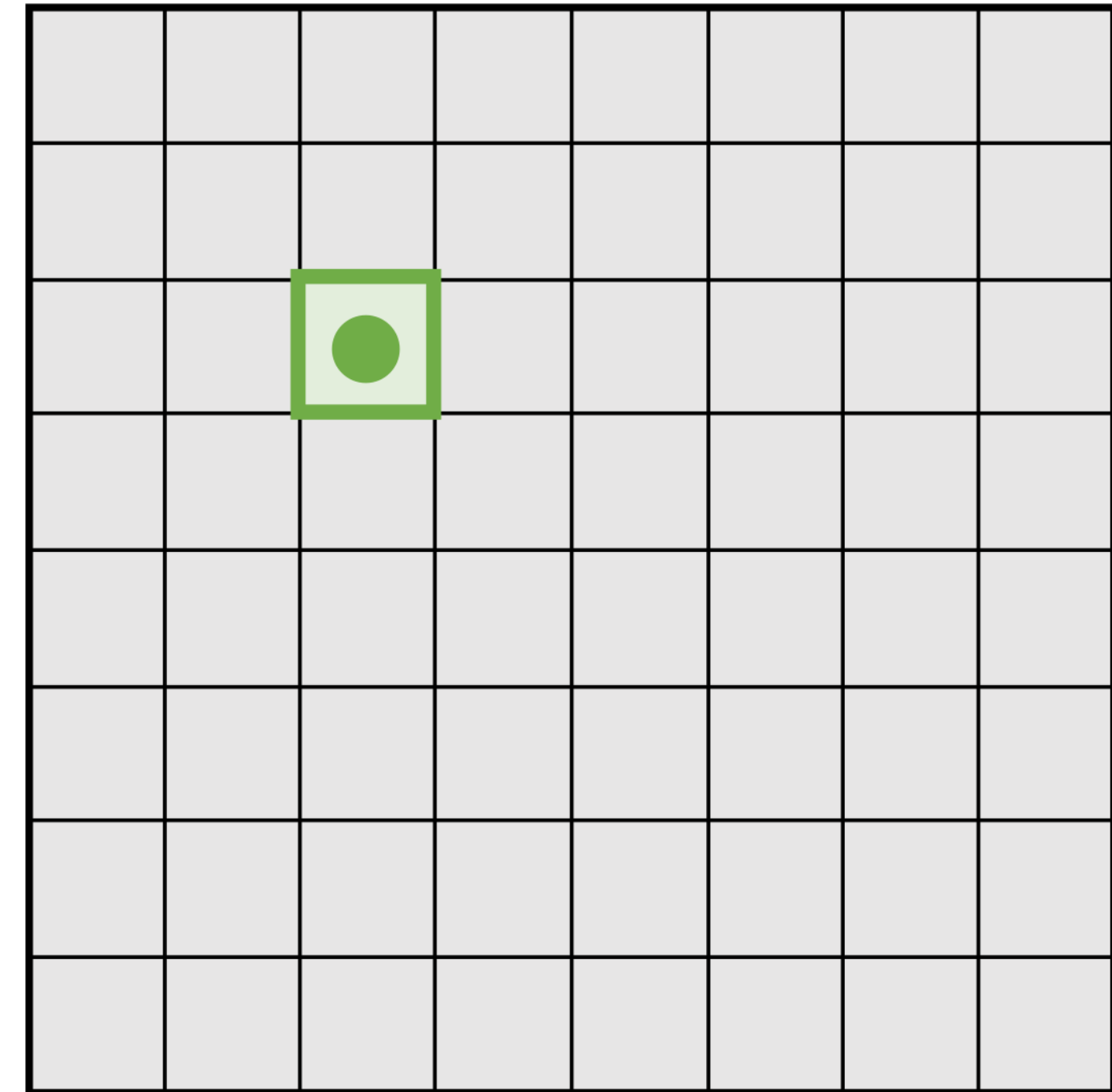


Input Image: 8 x 8

Every position in the output feature map depends on a 5x5 receptive field in the input

3x3 Conv
Stride 1, pad 1

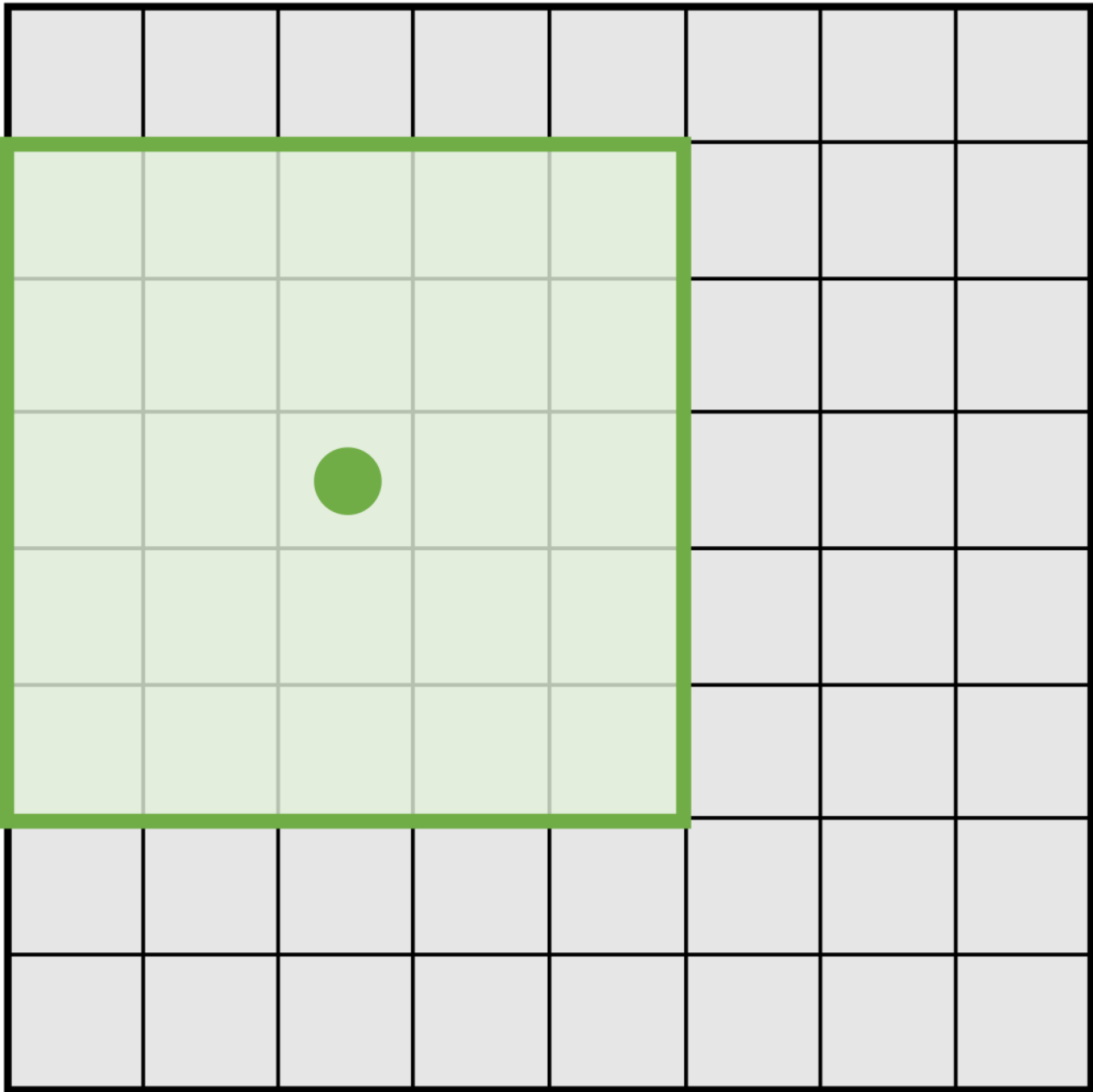
3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8



Recall: Receptive Fields

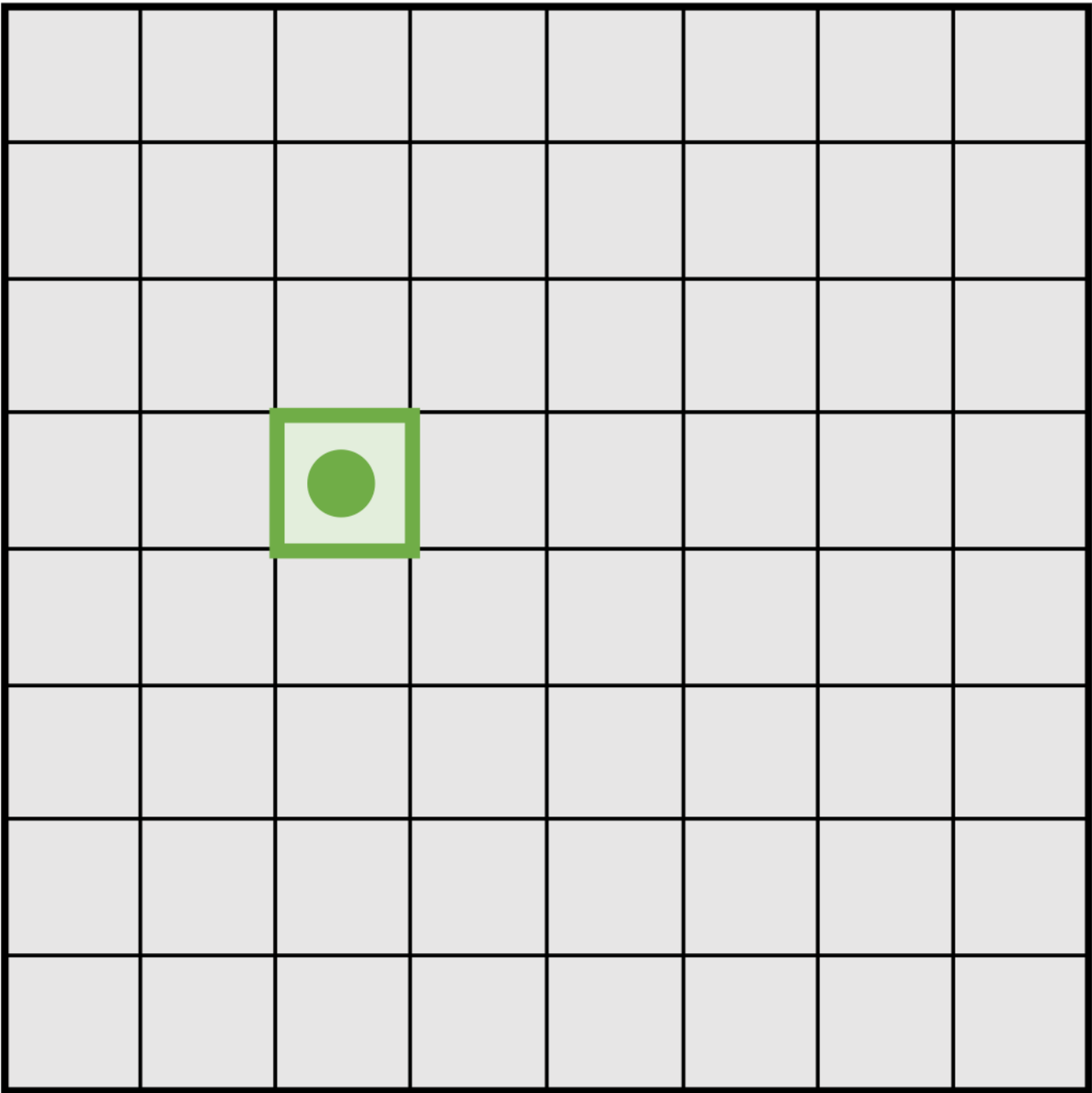


Input Image: 8 x 8

Moving one unit in the output space also moves the receptive field by one

3x3 Conv
Stride 1, pad 1

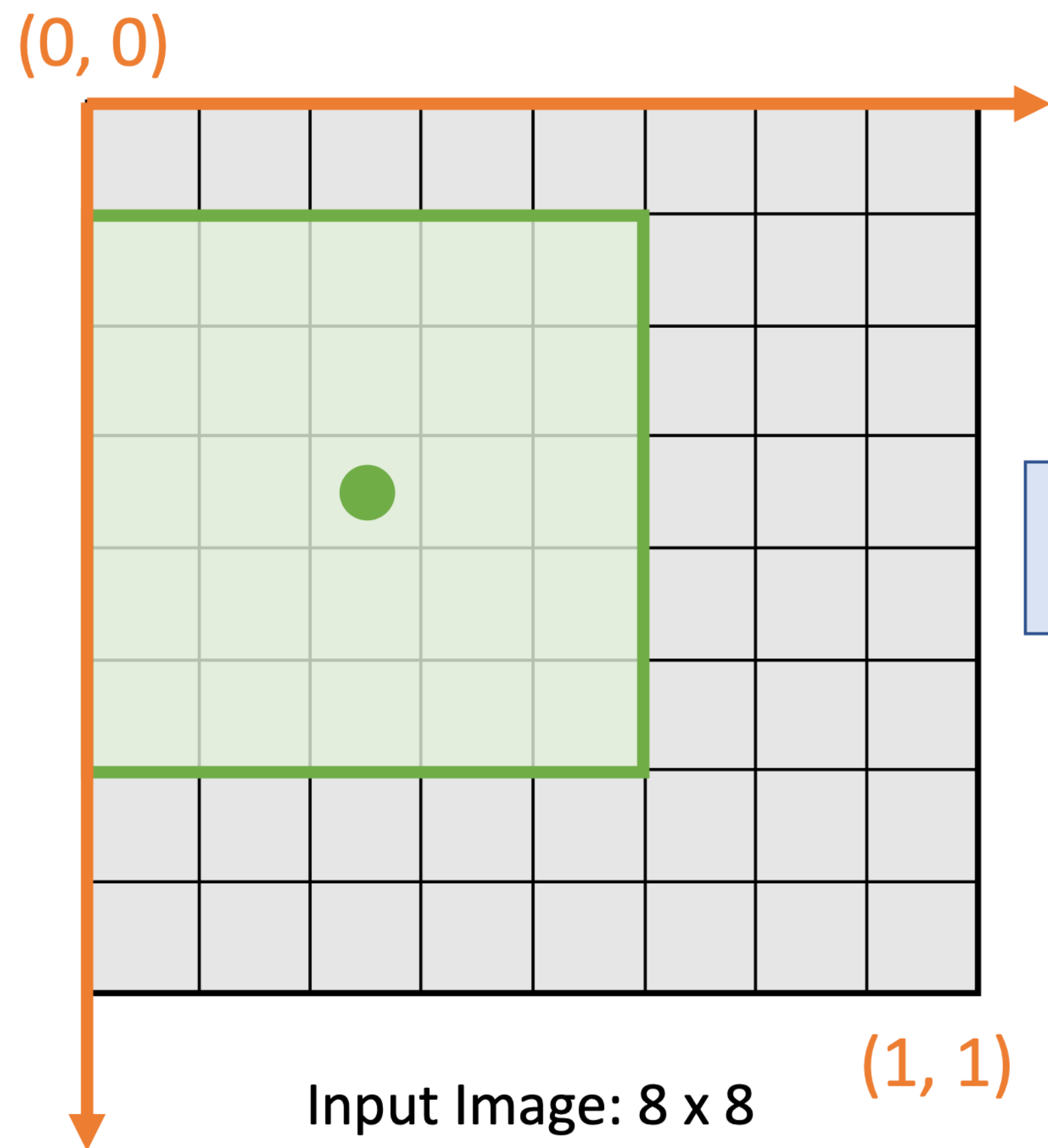
3x3 Conv
Stride 1, pad 1



Output Image: 8 x 8



Recall: Receptive Fields

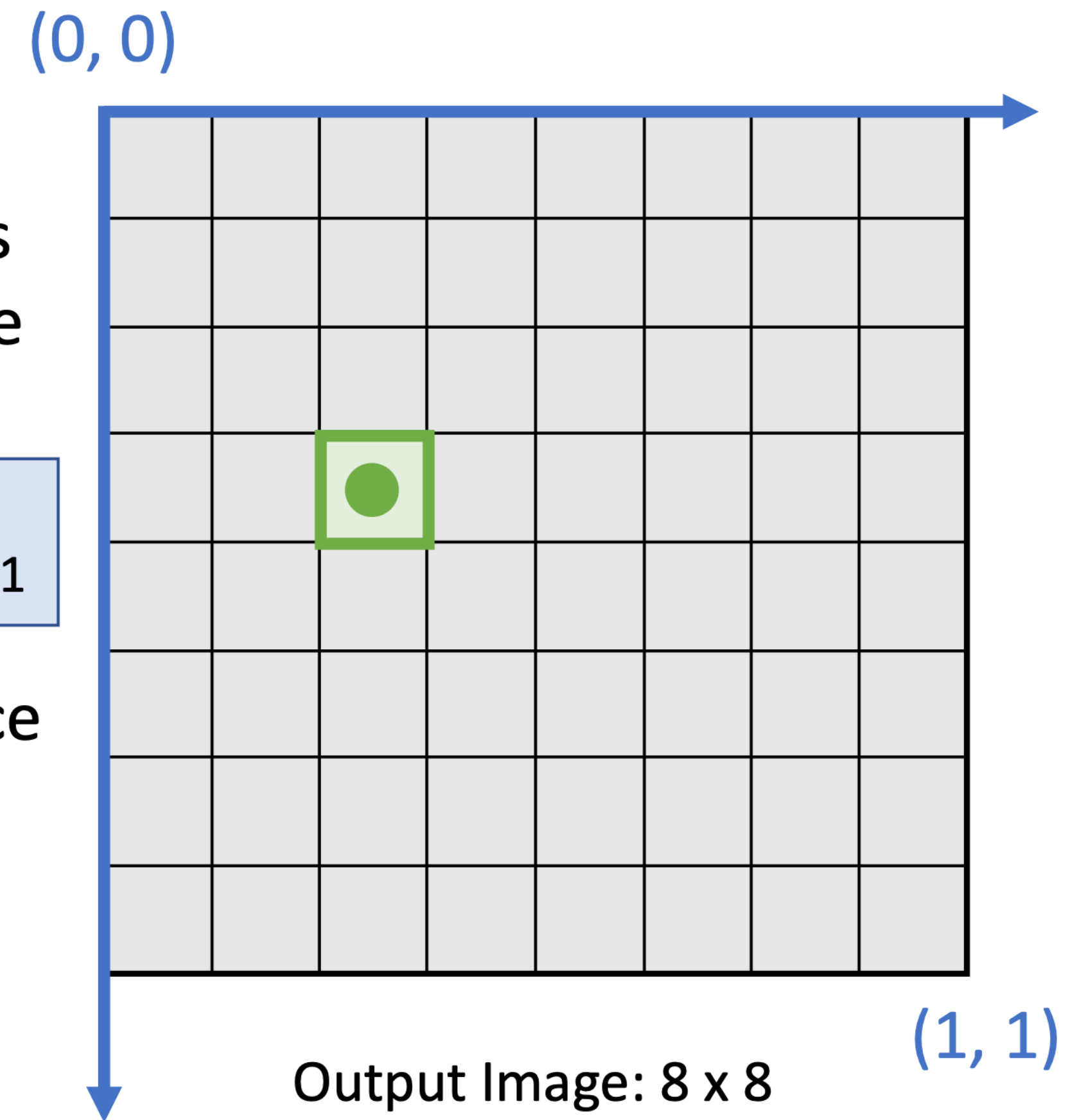


Moving one unit in the output space also moves the receptive field by one

3x3 Conv
Stride 1, pad 1

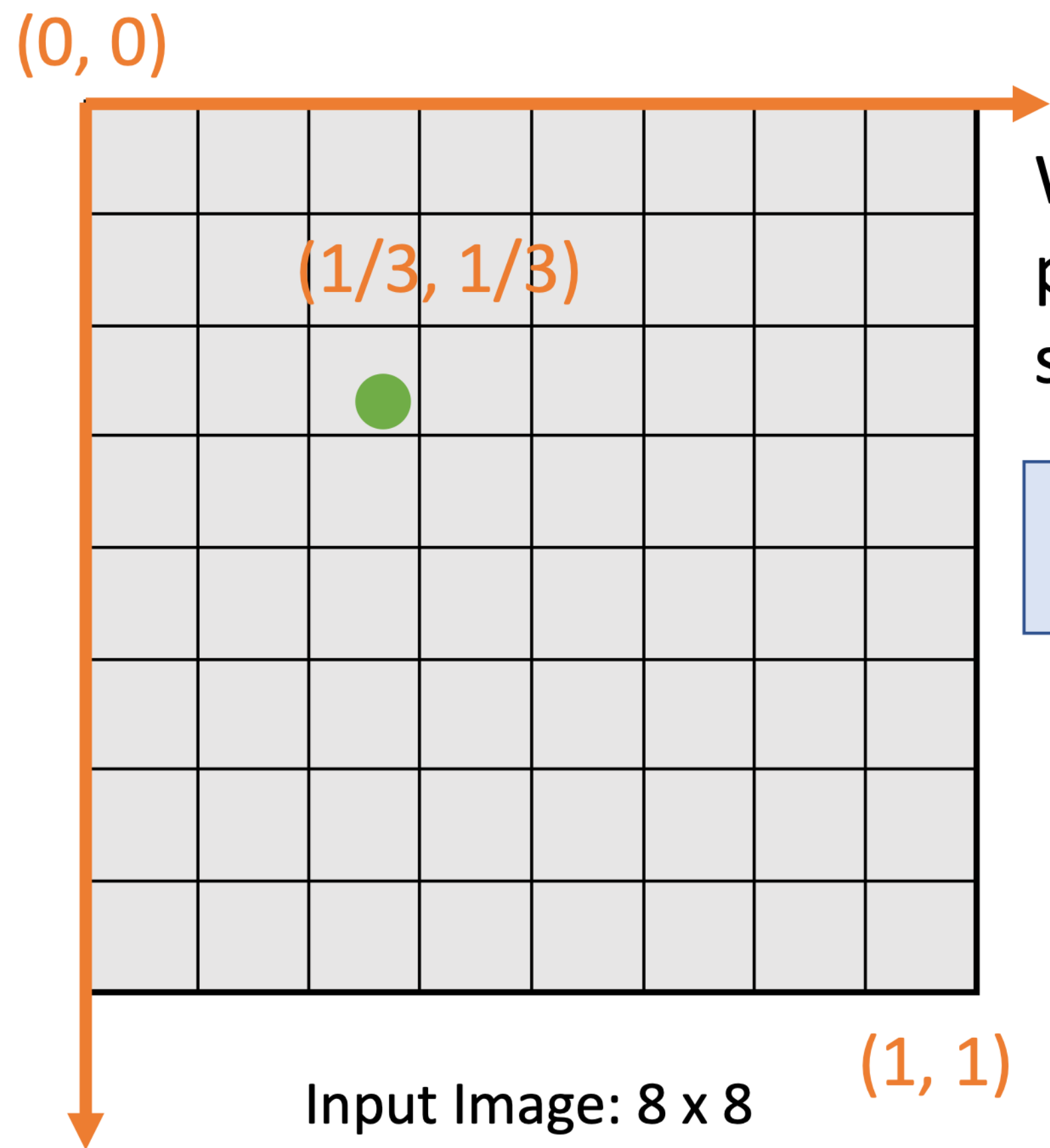
3x3 Conv
Stride 1, pad 1

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**





Projecting Points

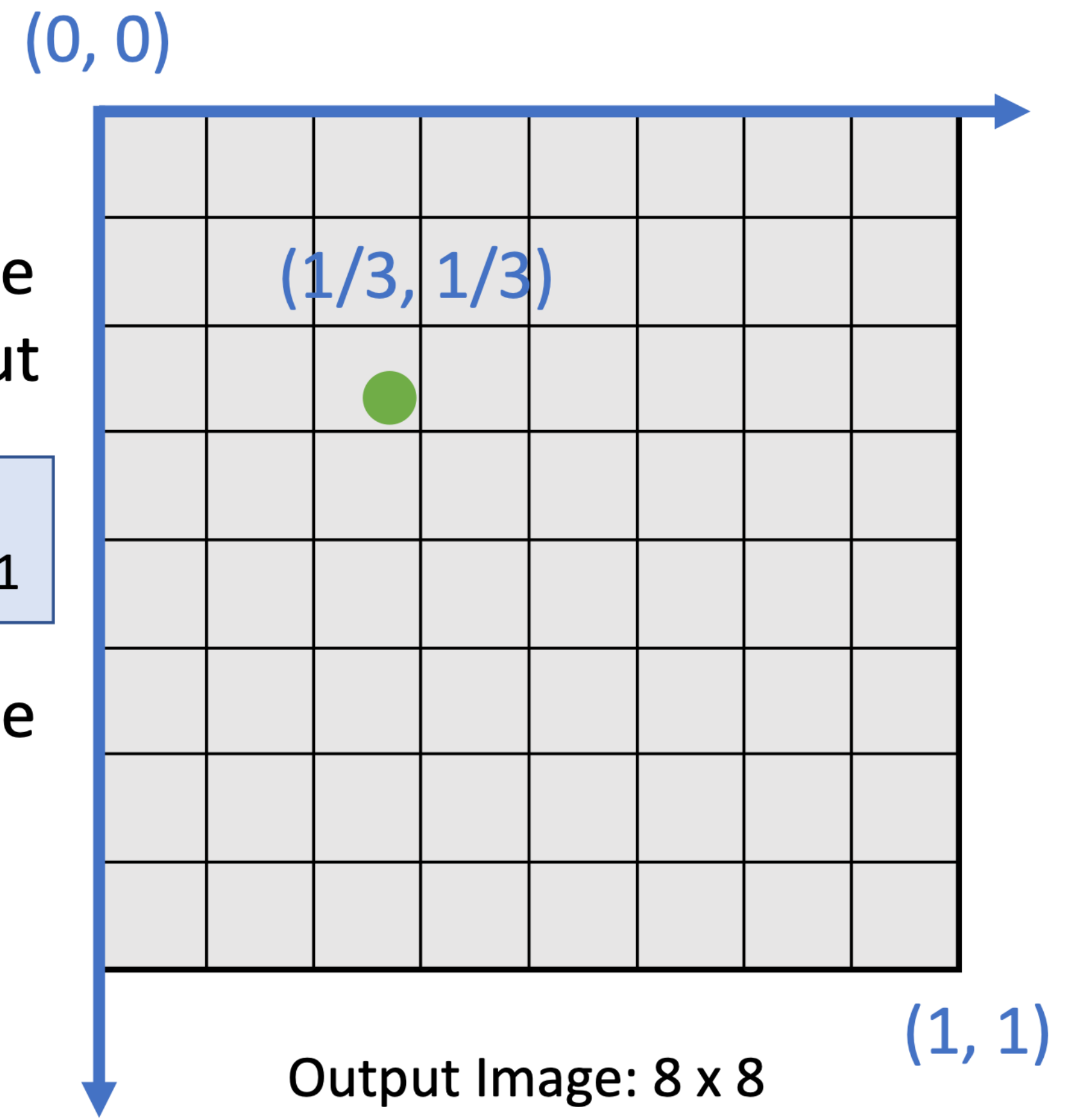


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

3x3 Conv
Stride 1, pad 1

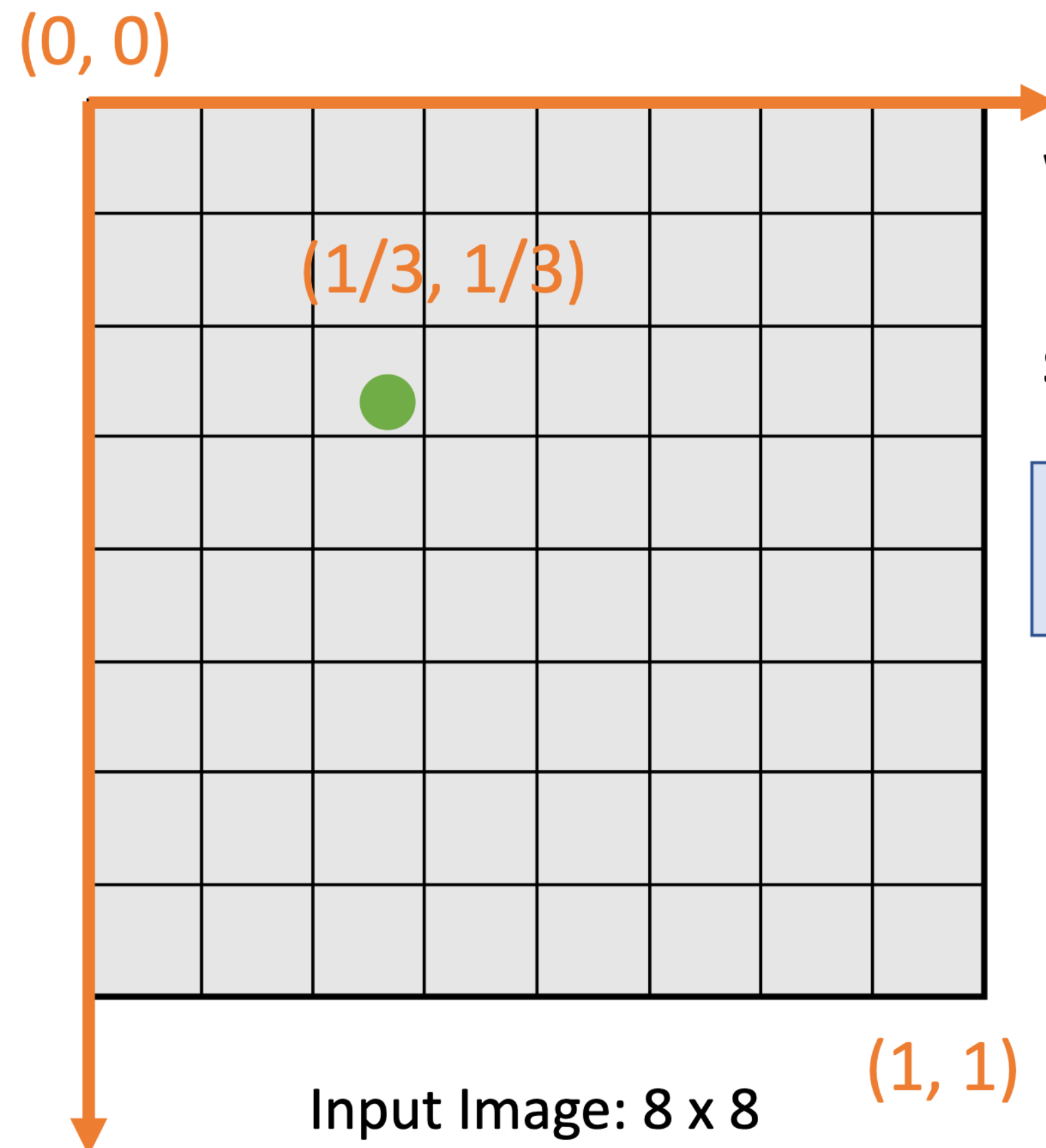
There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**





Projecting Points

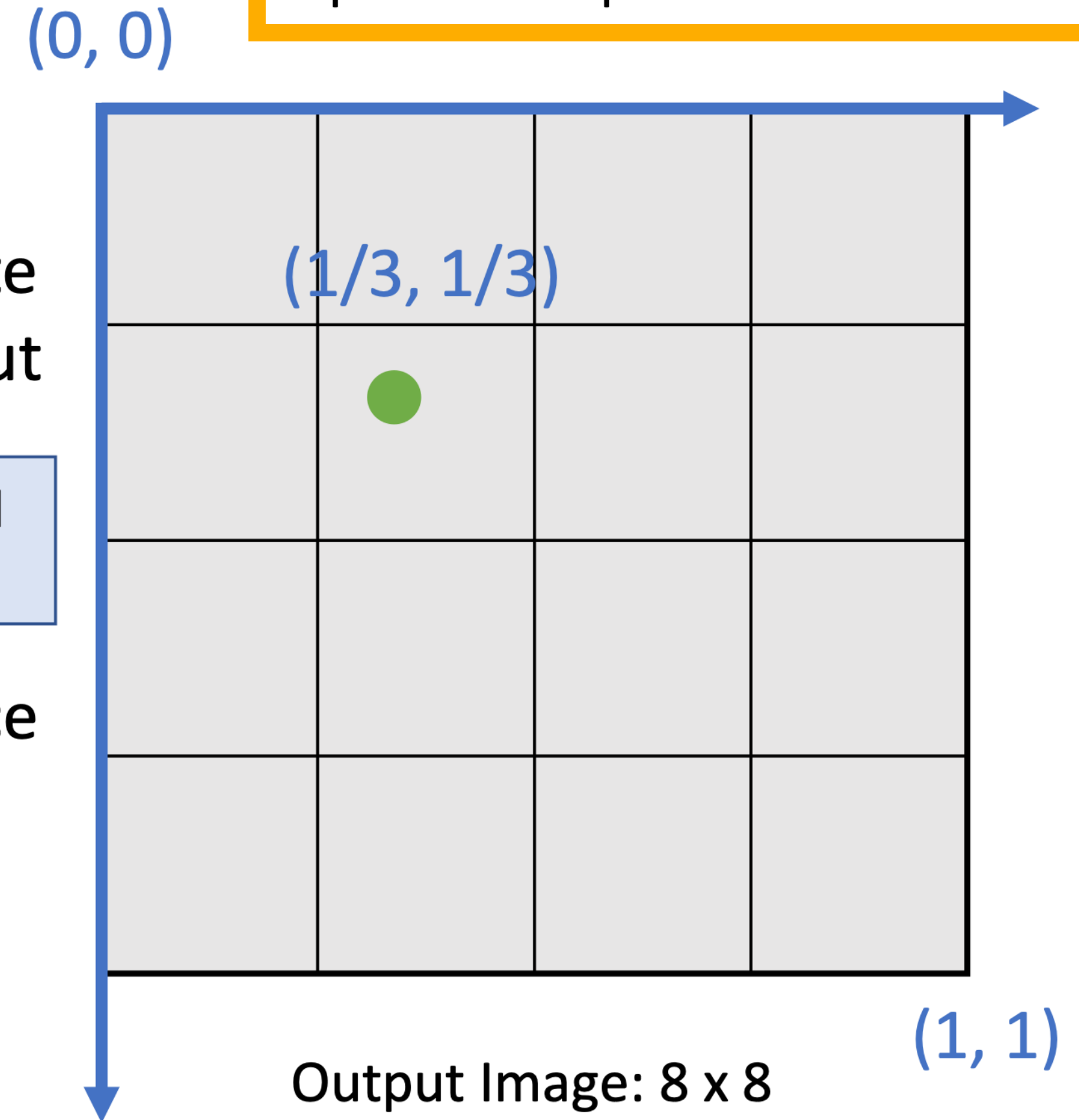
Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different



We can align arbitrary points between coordinate system of input and output

3x3 Conv Stride 1, pad 1 2x2 MaxPool Stride 2

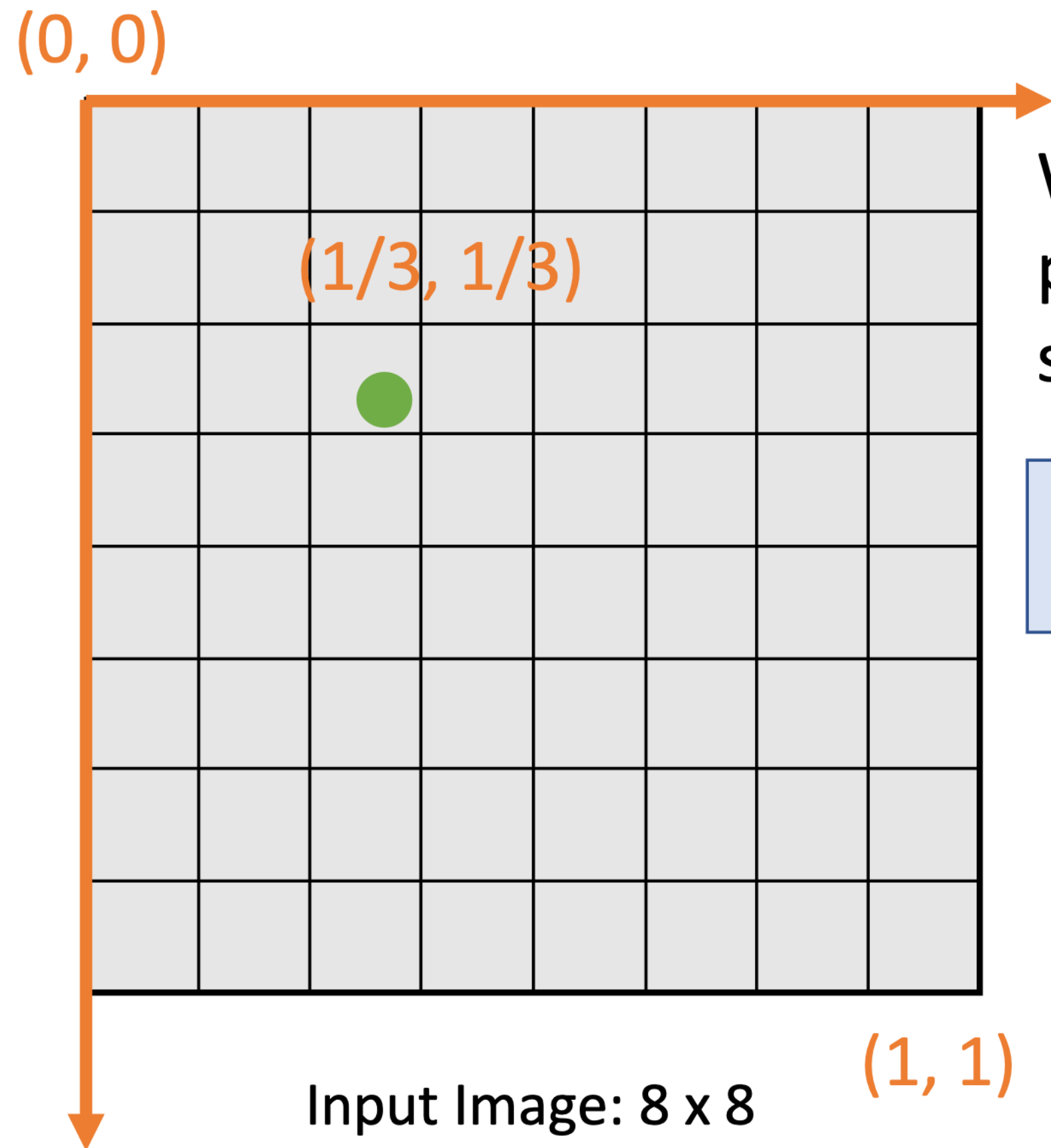
There is a correspondence between the coordinate system of the input and the coordinate system of the output





Projecting Points

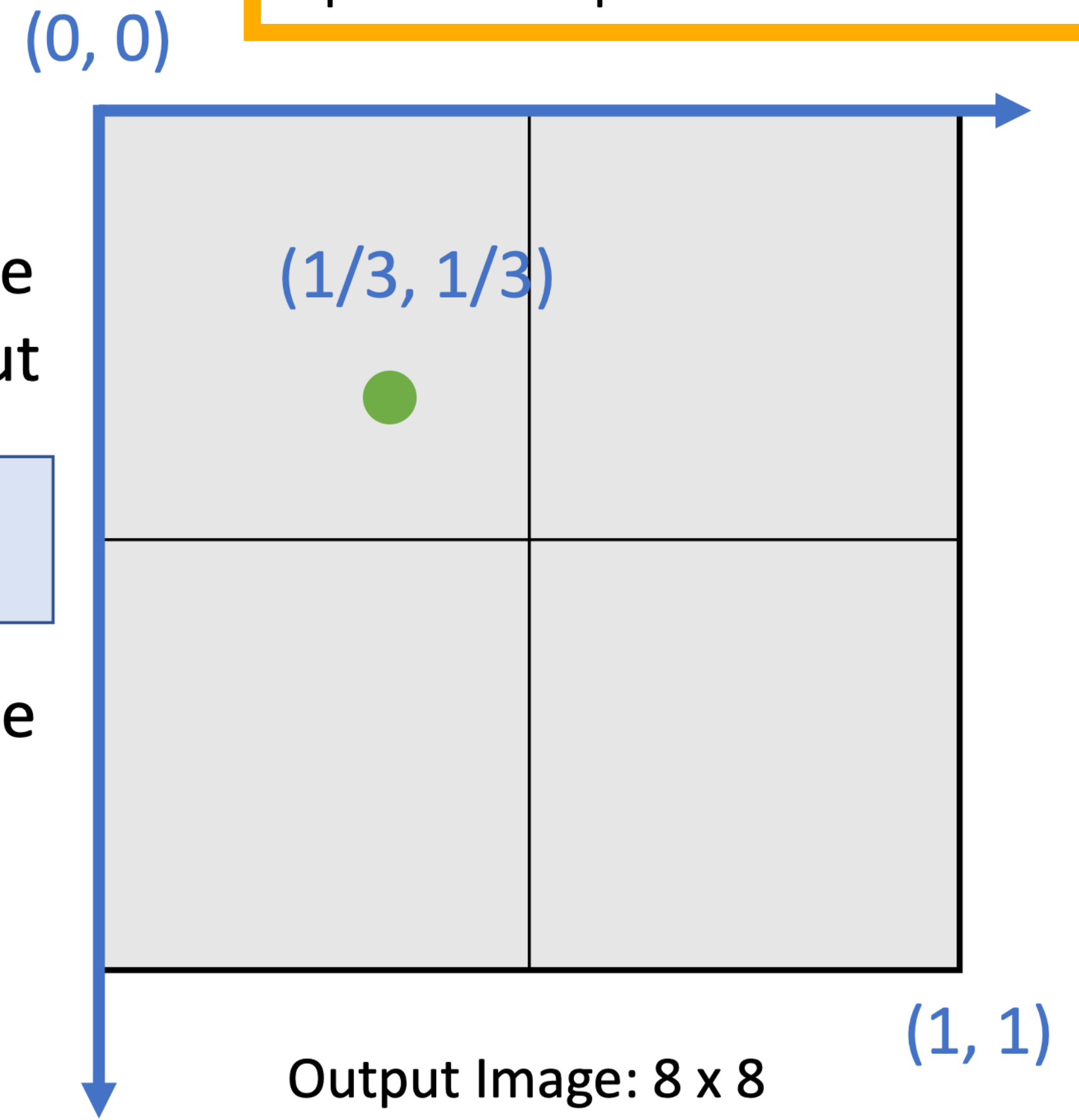
Same logic holds for more complicated CNNs, even if spatial resolution of input and output are different



We can align arbitrary points between coordinate system of input and output

- 3x3 Conv
Stride 1, pad 1
- 4x4 MaxPool
Stride 4

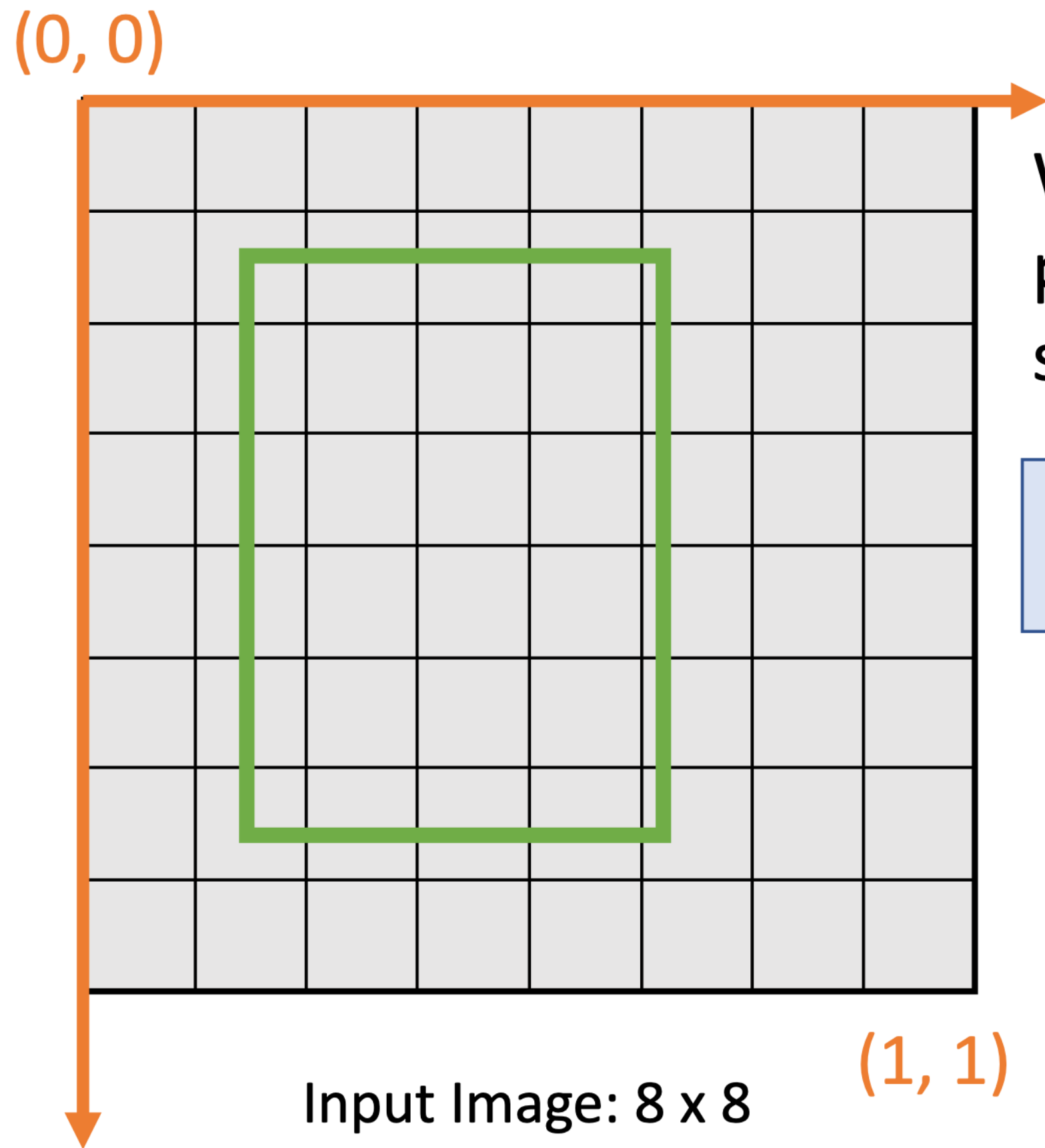
There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**





Projecting Points

We can use this idea to project **bounding boxes** between an input image and a feature map

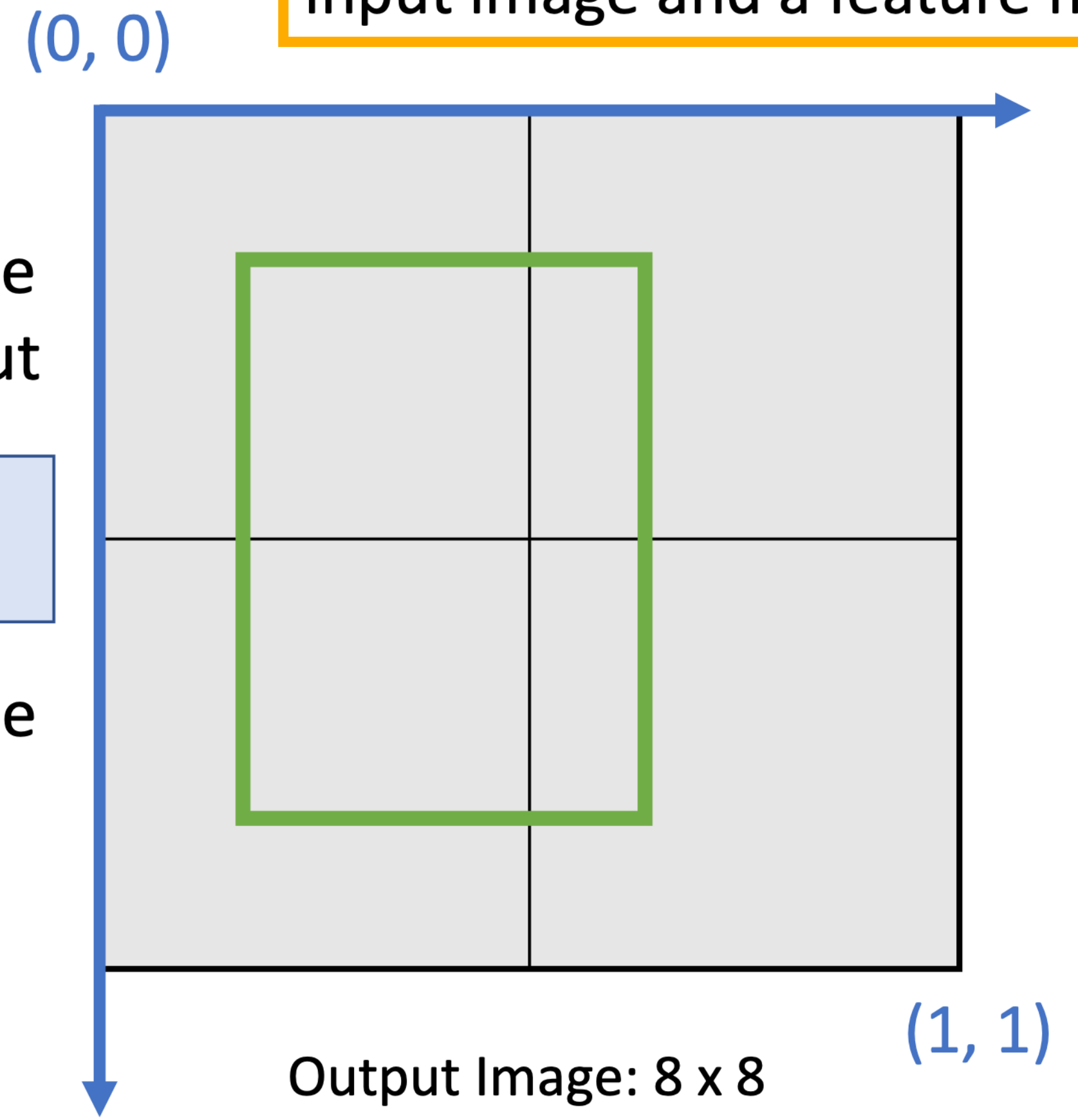


We can align arbitrary points between coordinate system of input and output

3x3 Conv
Stride 1, pad 1

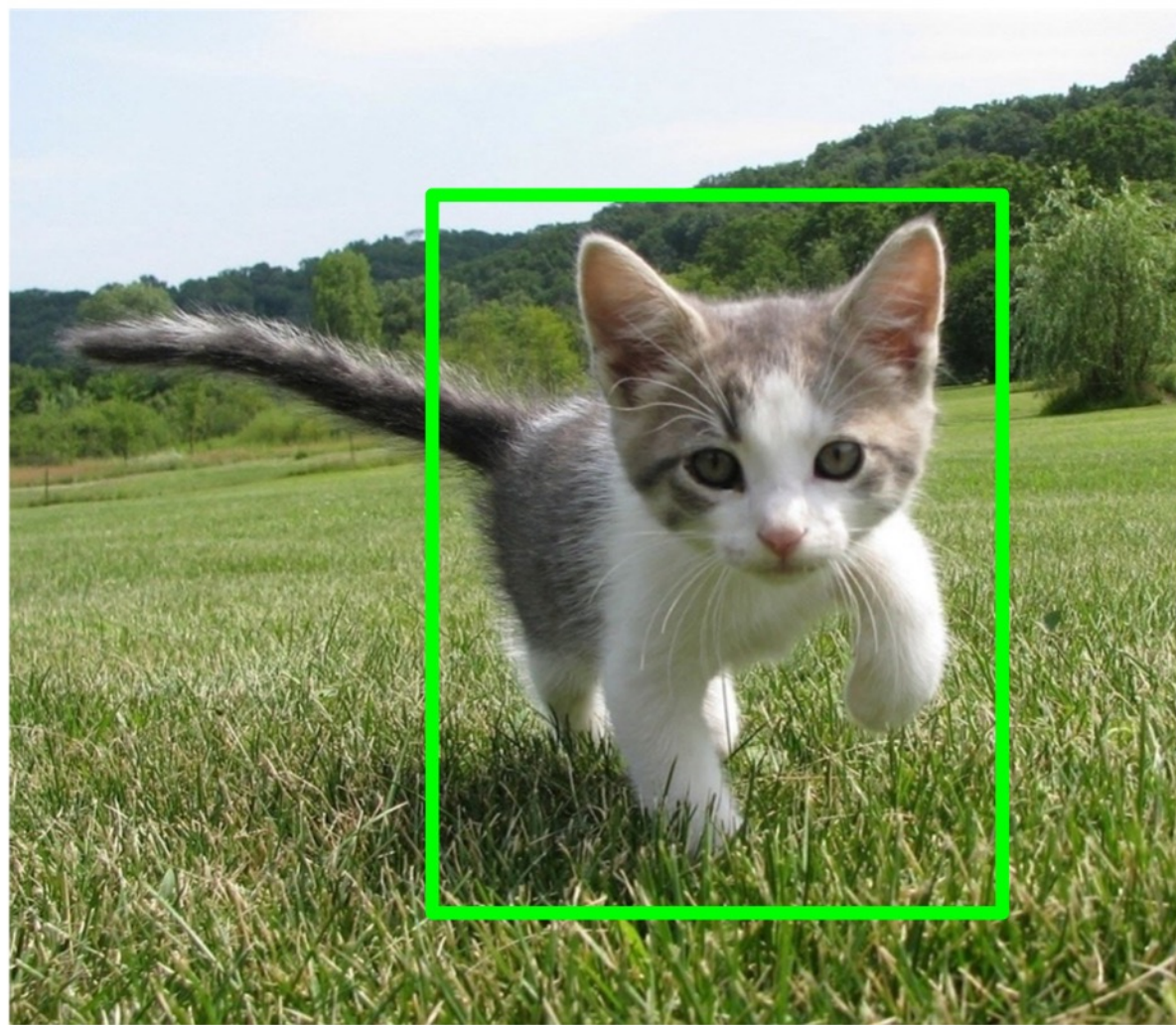
4x4 MaxPool
Stride 4

There is a correspondence between the **coordinate system of the input** and the **coordinate system of the output**





Cropping Features: RoI Pool



Input Image
(e.g. 3 x 640 x 480)

CNN

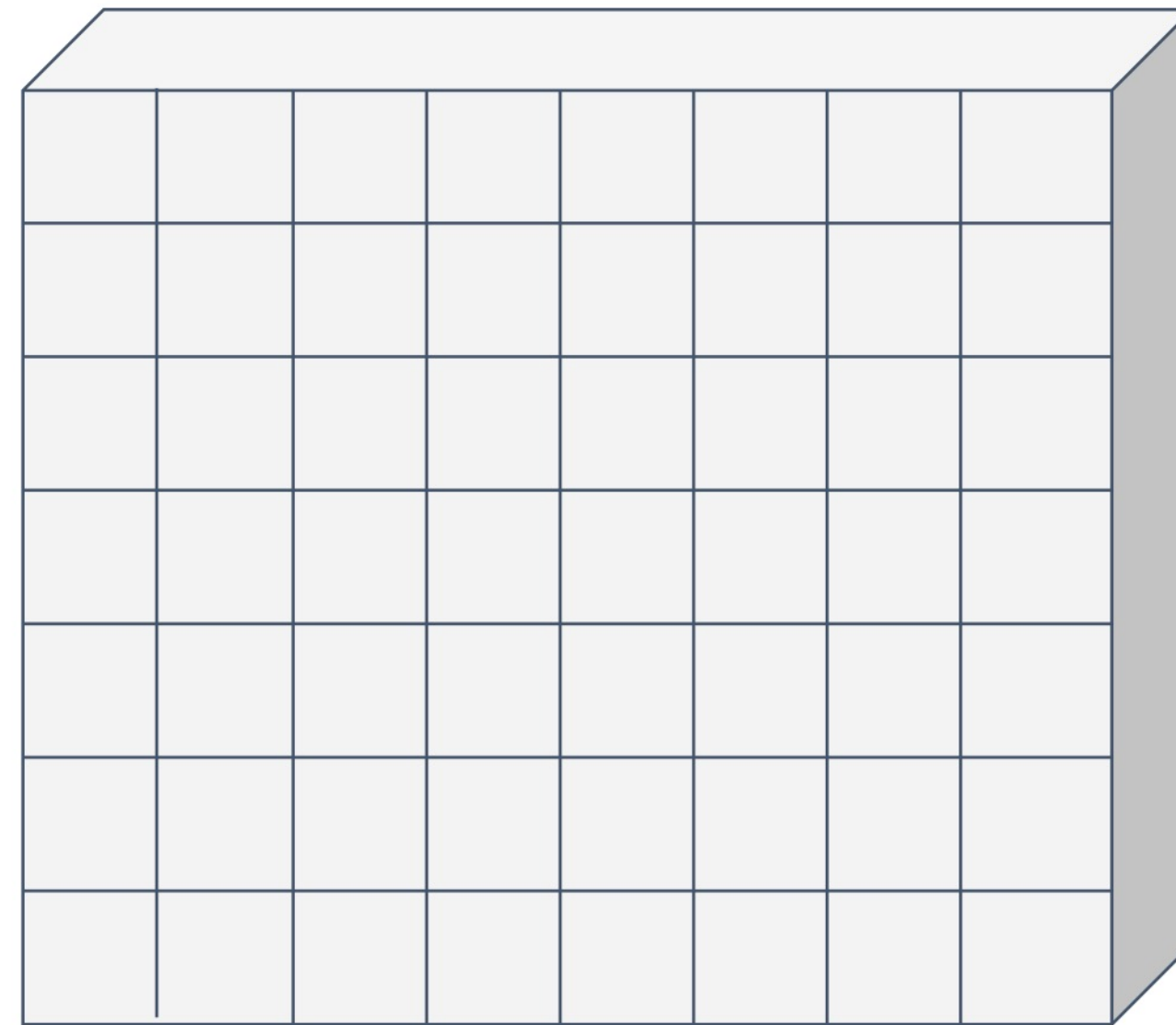
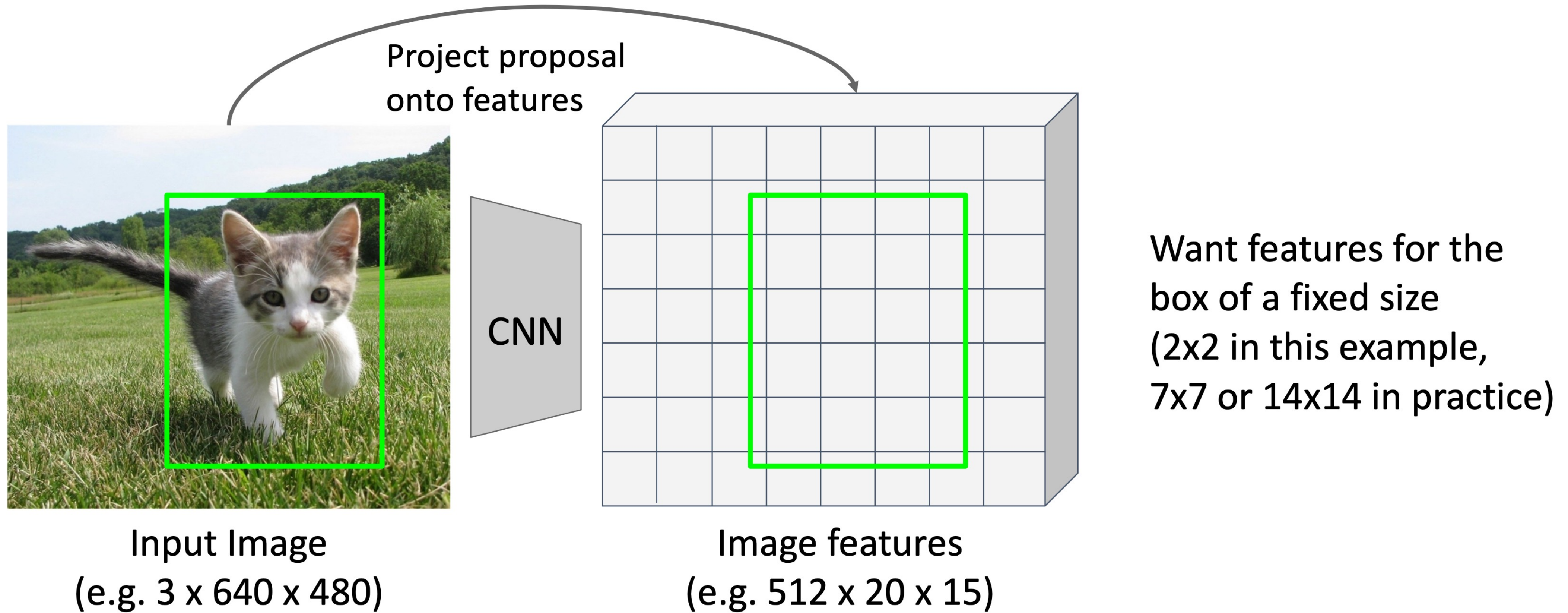


Image features
(e.g. 512 x 20 x 15)

Want features for the box of a fixed size (2x2 in this example, 7x7 or 14x14 in practice)

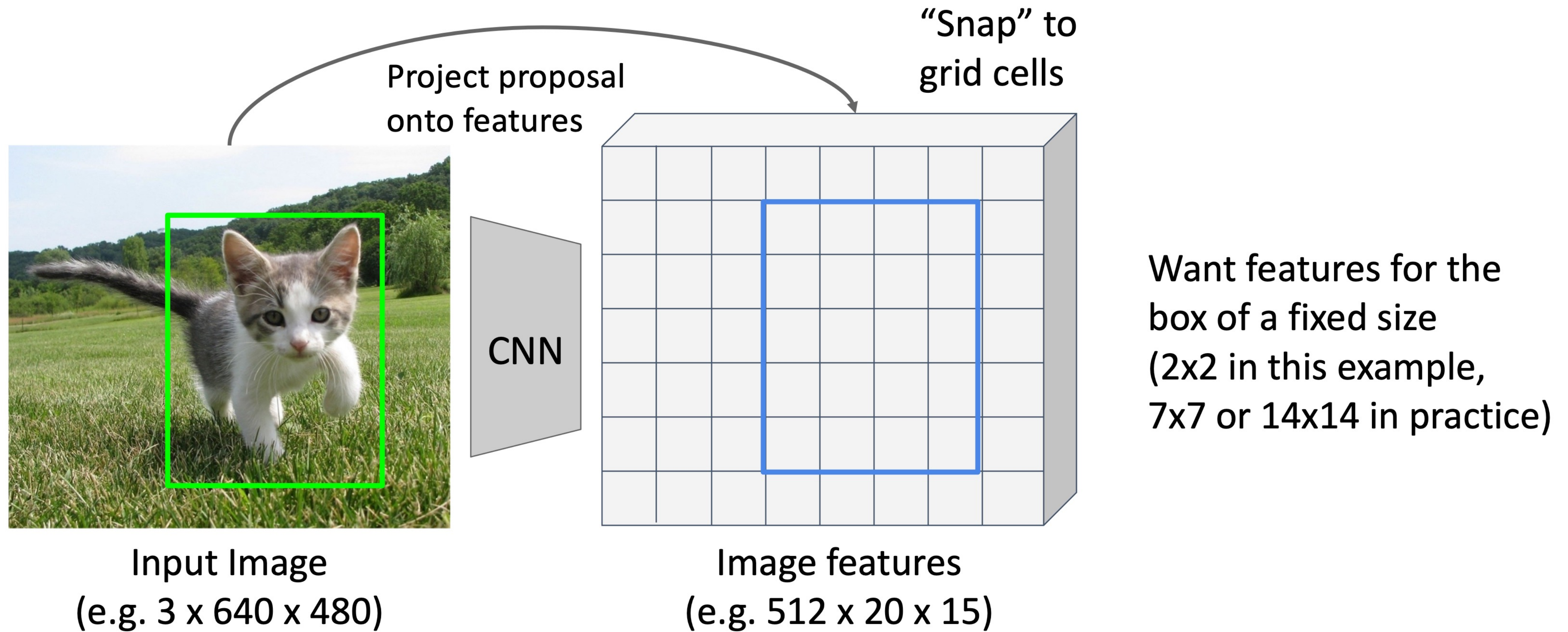


Cropping Features: RoI Pool



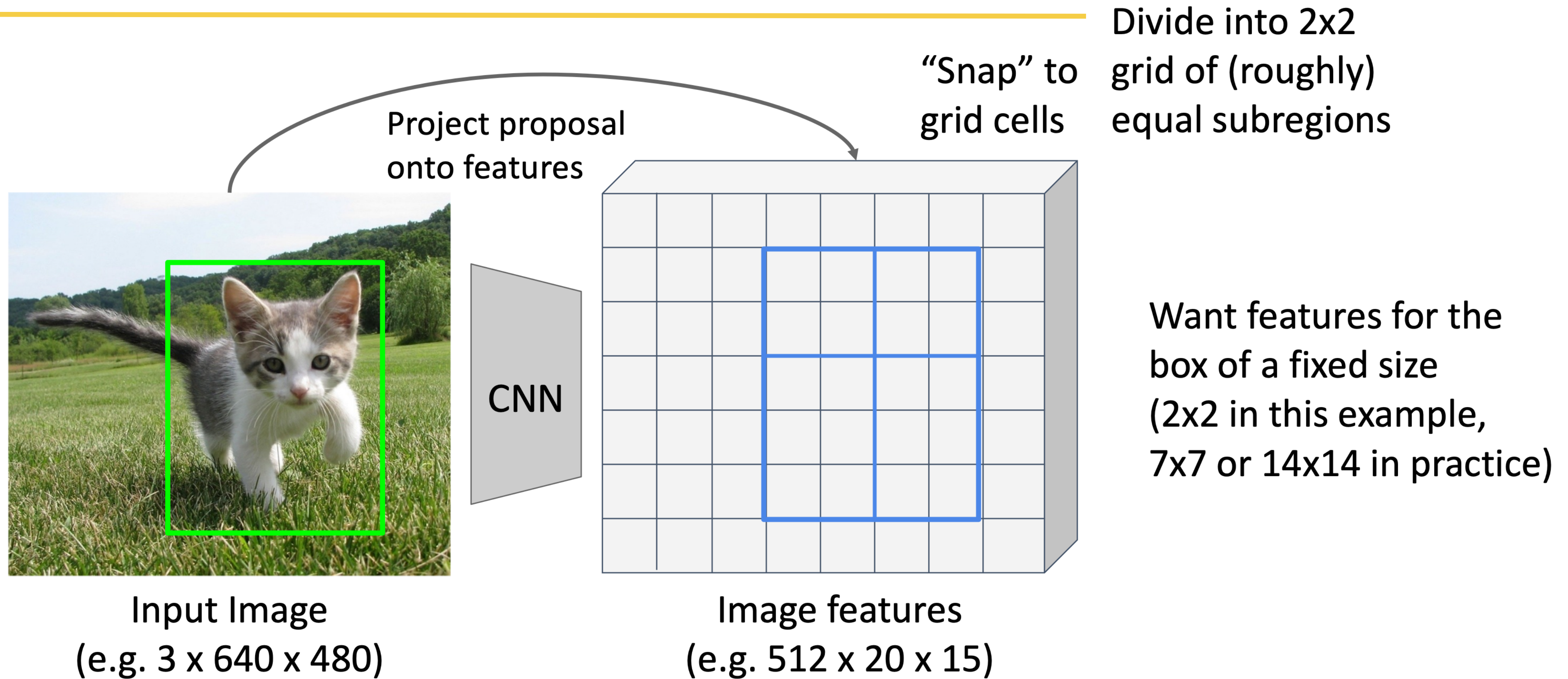


Cropping Features: RoI Pool



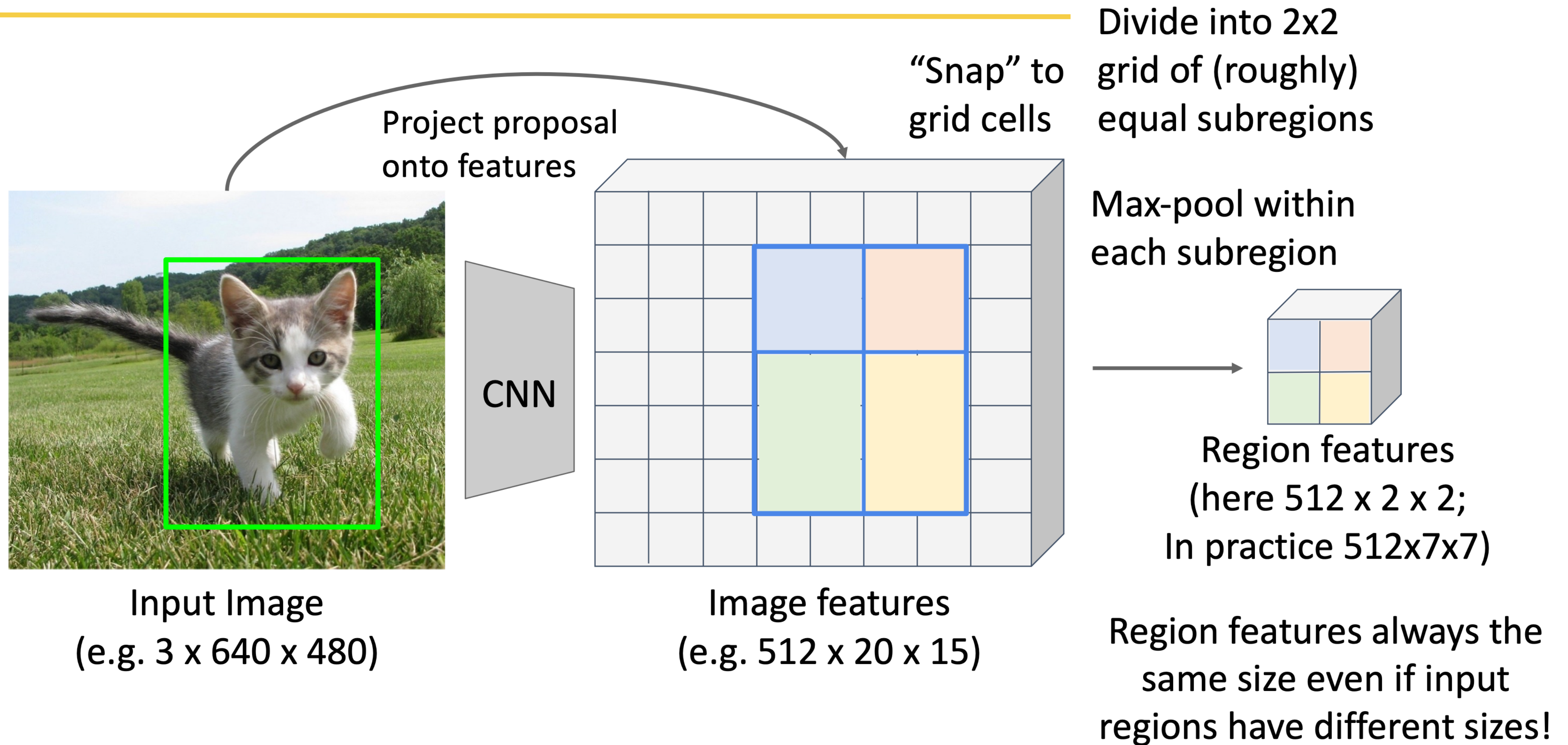


Cropping Features: RoI Pool





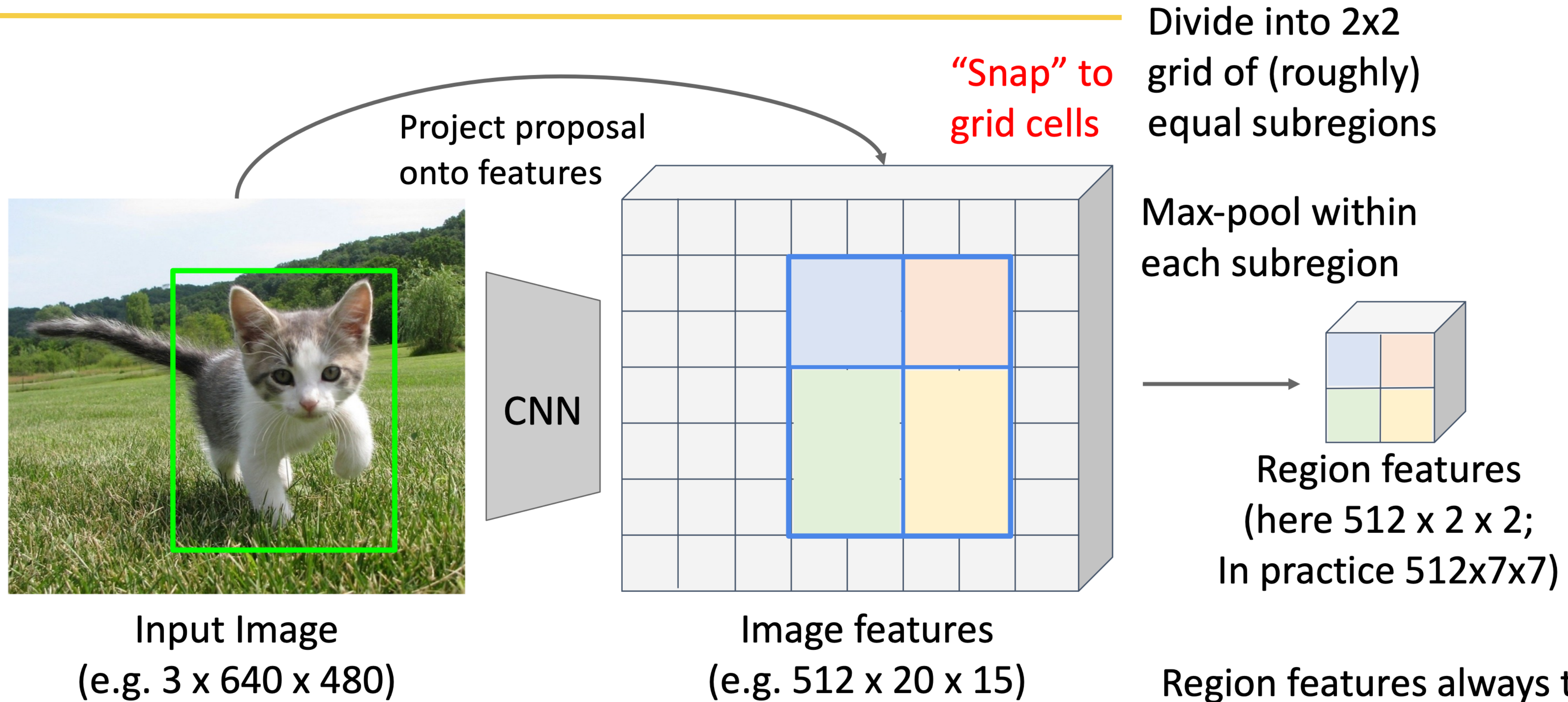
Cropping Features: RoI Pool



Girshick, "Fast R-CNN", ICCV 2015.



Cropping Features: RoI Pool



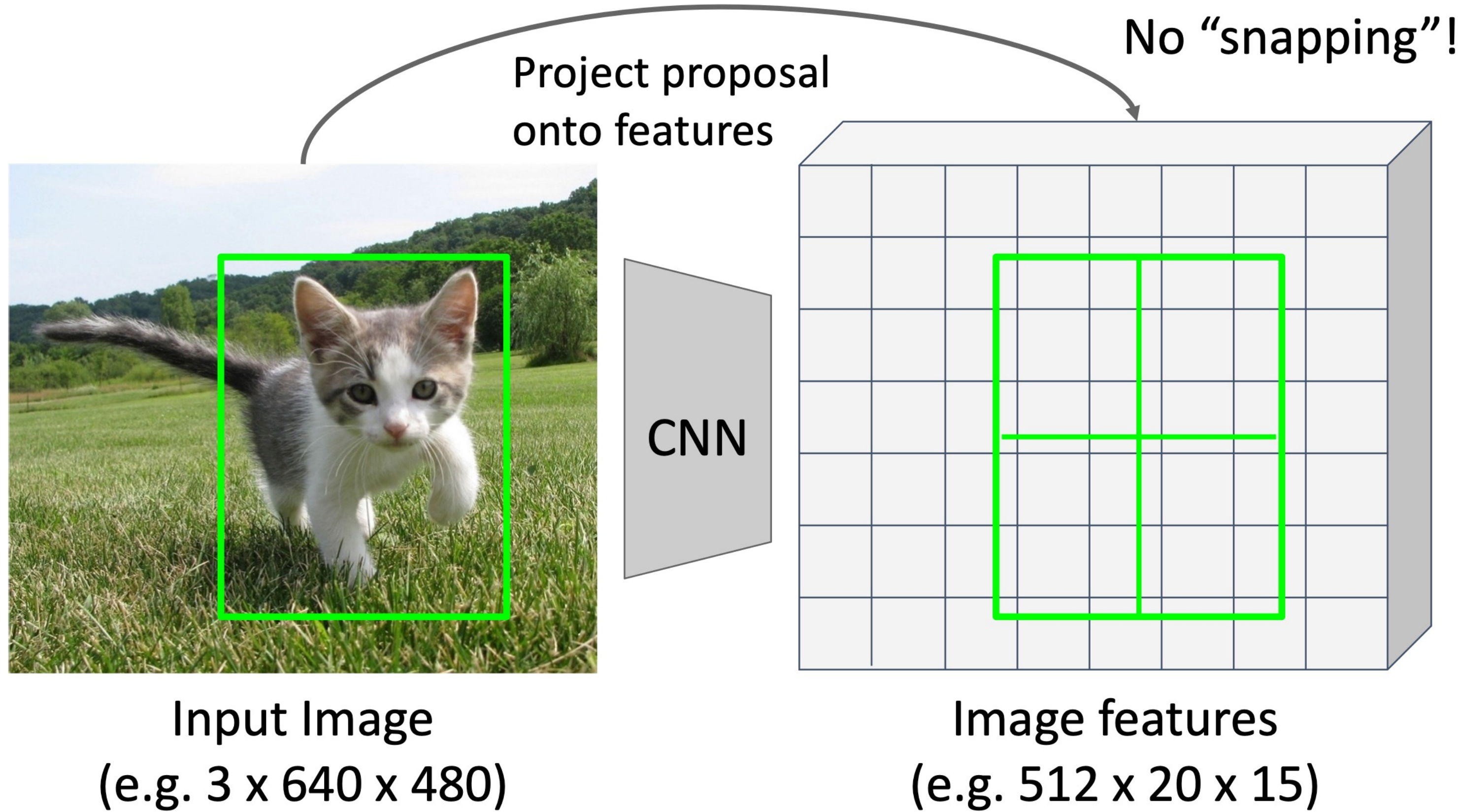
Problem: Slight misalignment due to snapping; different-sized subregions is weird

Girshick, “Fast R-CNN”, ICCV 2015.



Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)

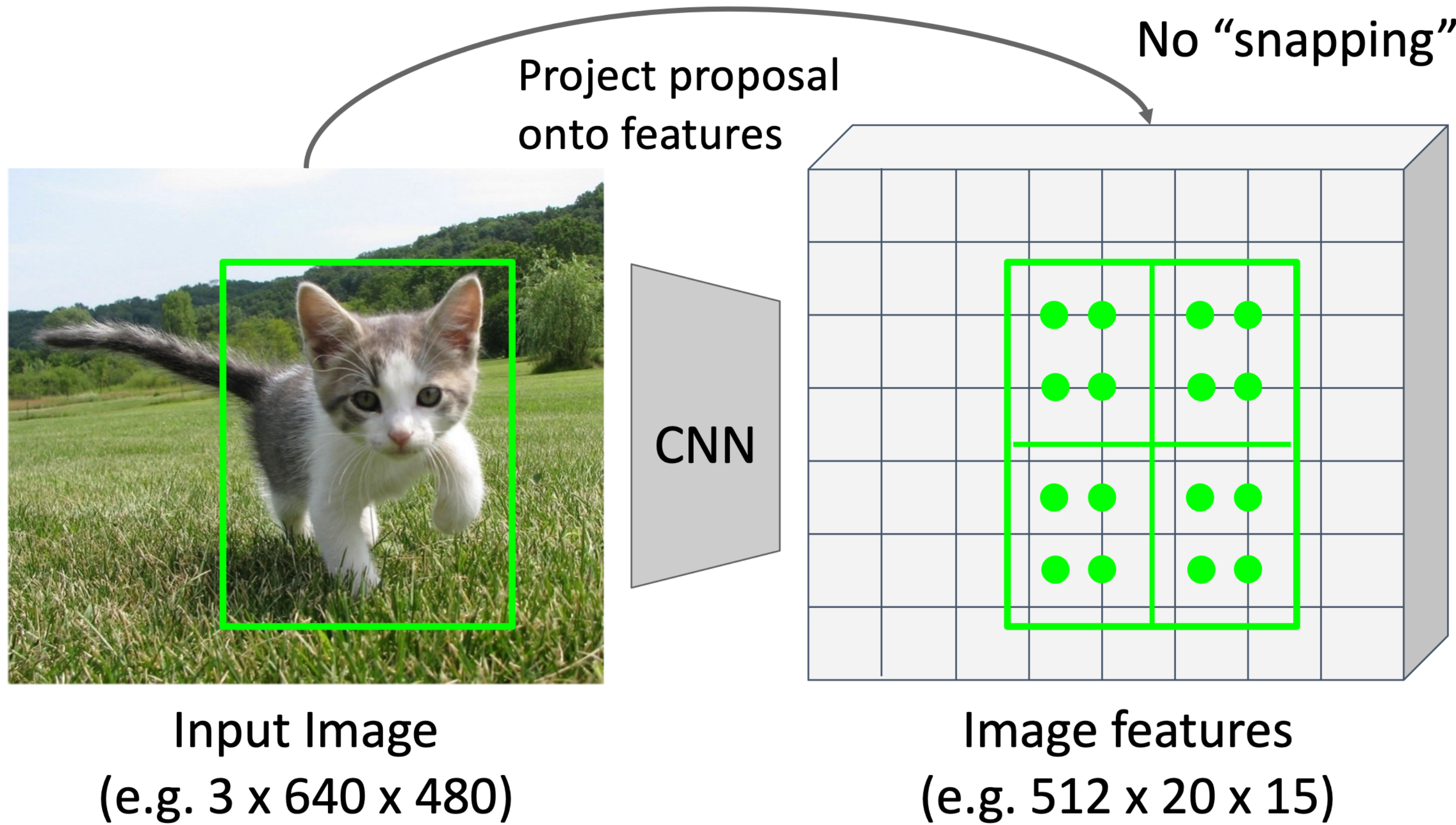


Want features for the box of a fixed size
(2x2 in this example,
7x7 or 14x14 in practice)



Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)

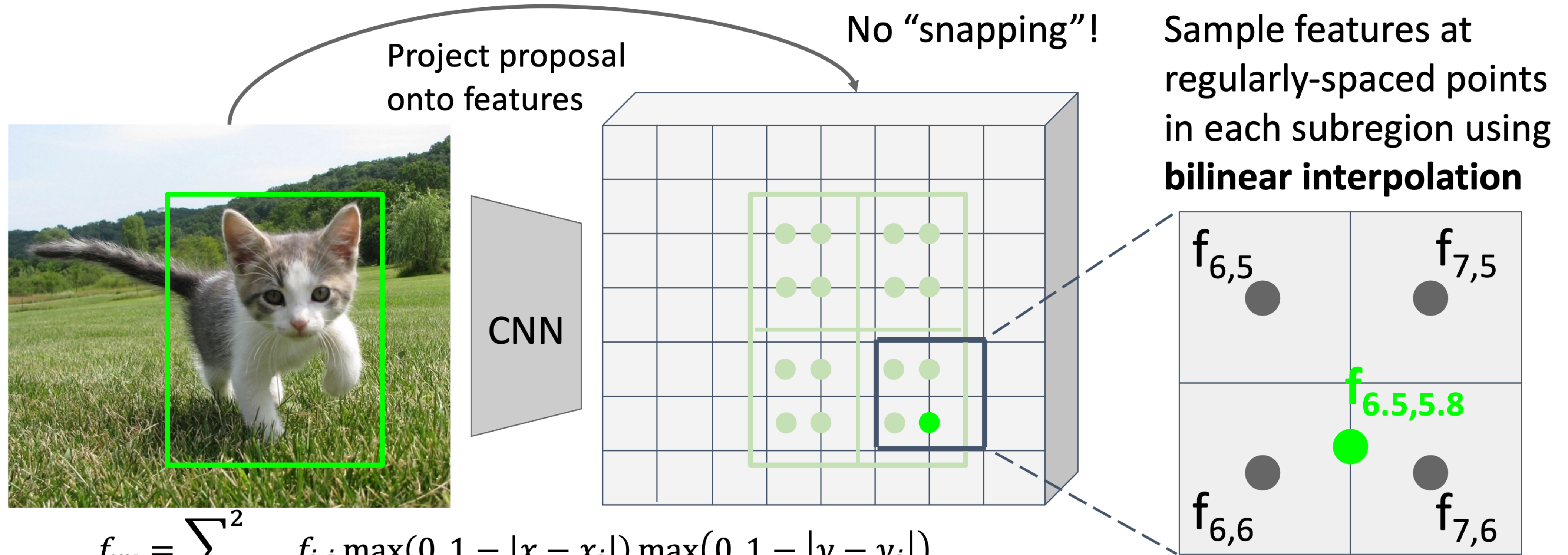


Sample features at regularly-spaced points in each subregion using **bilinear interpolation**



Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



$$f_{xy} = \sum_{i,j=1}^2 f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_j|)$$

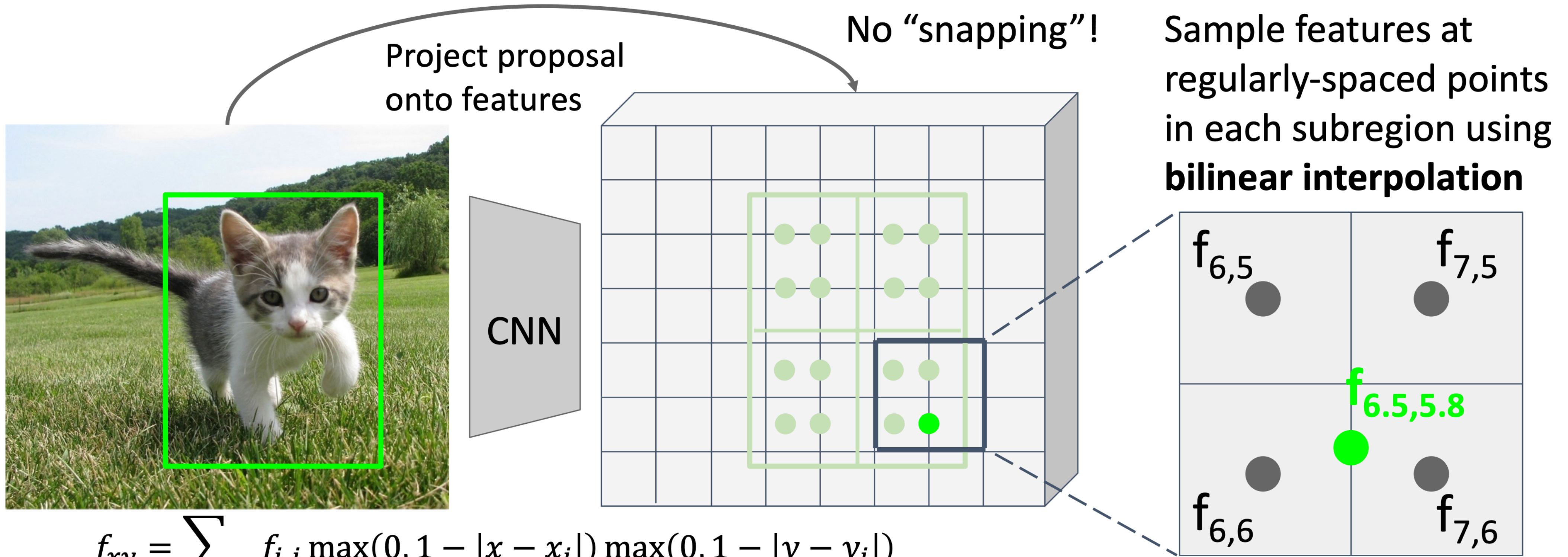
Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align

Divide into equal-sized subregions
(may not be aligned to grid!)



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

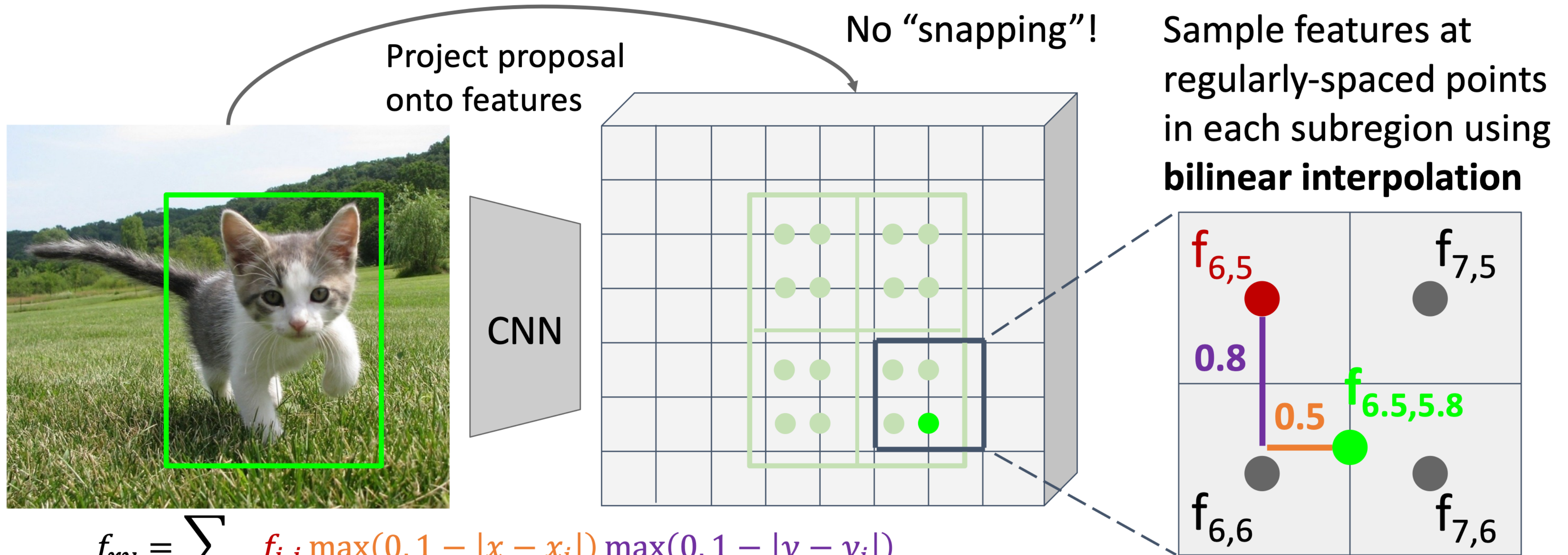
$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

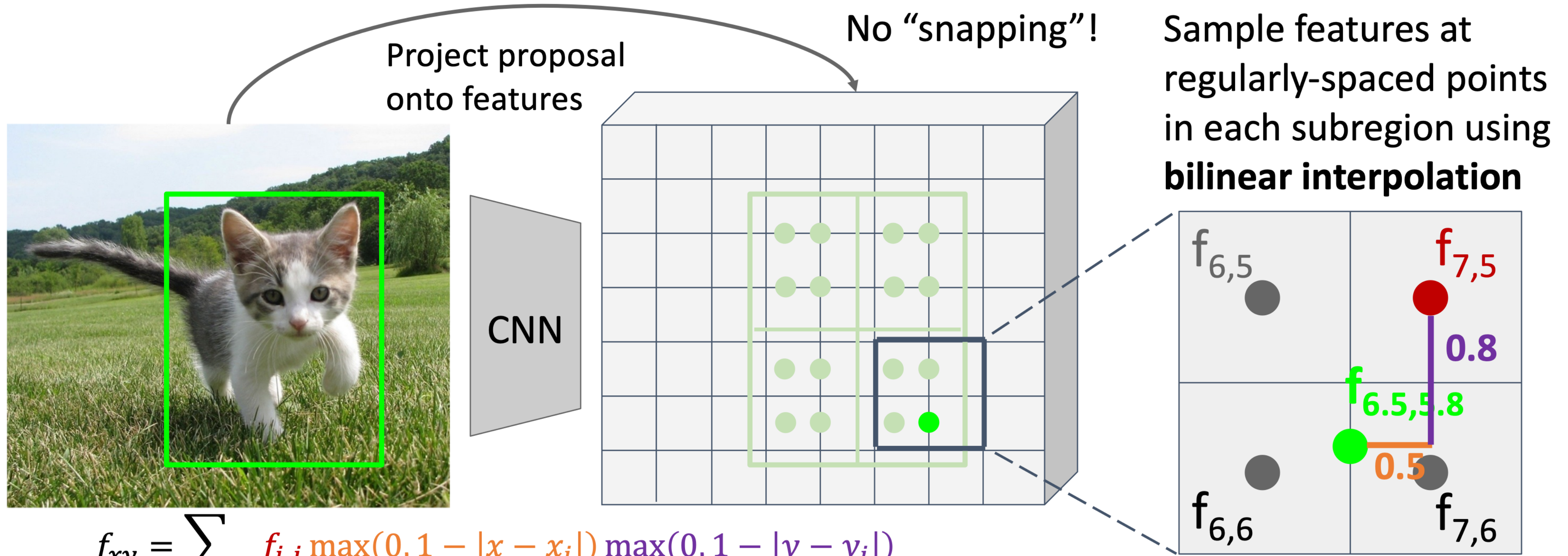
$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

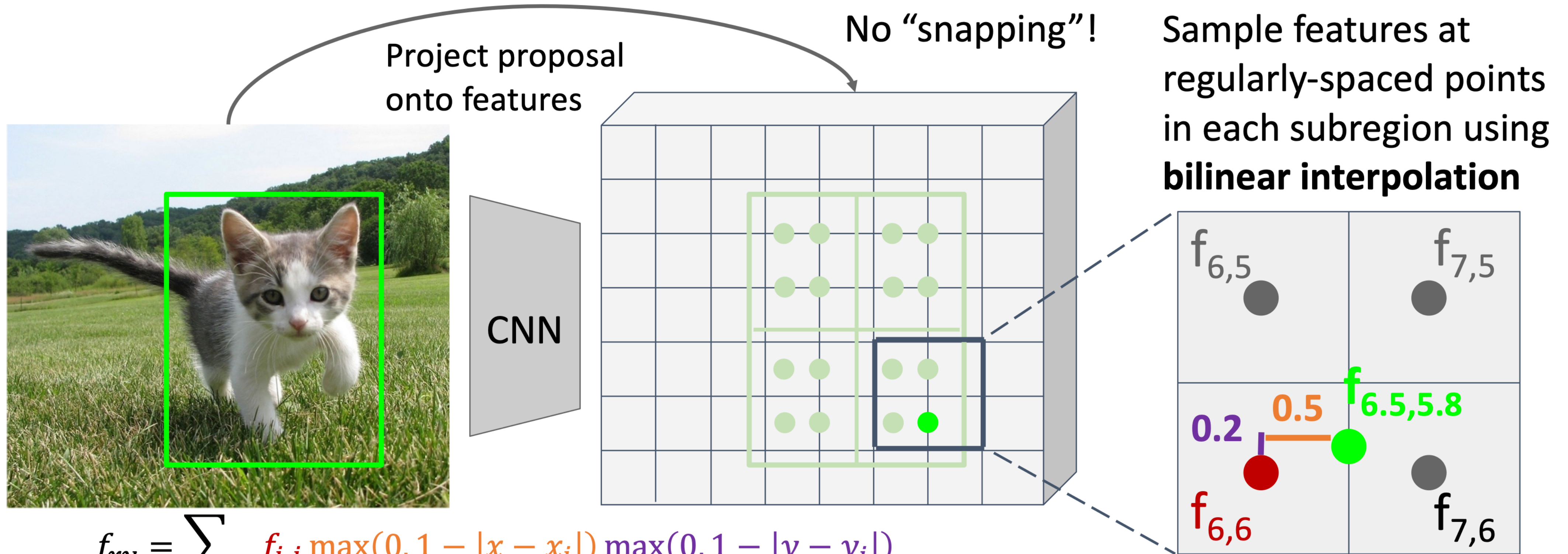
$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

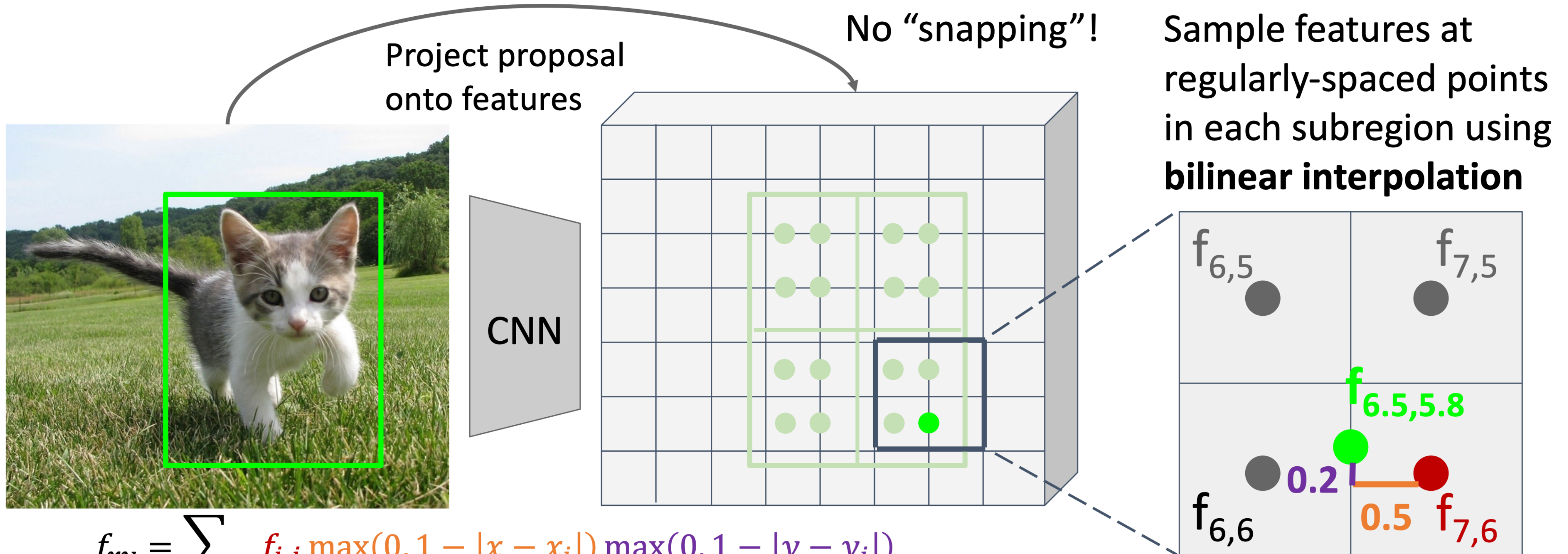
$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align



$$f_{xy} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - x_i|) \max(0, 1 - |y - y_i|)$$

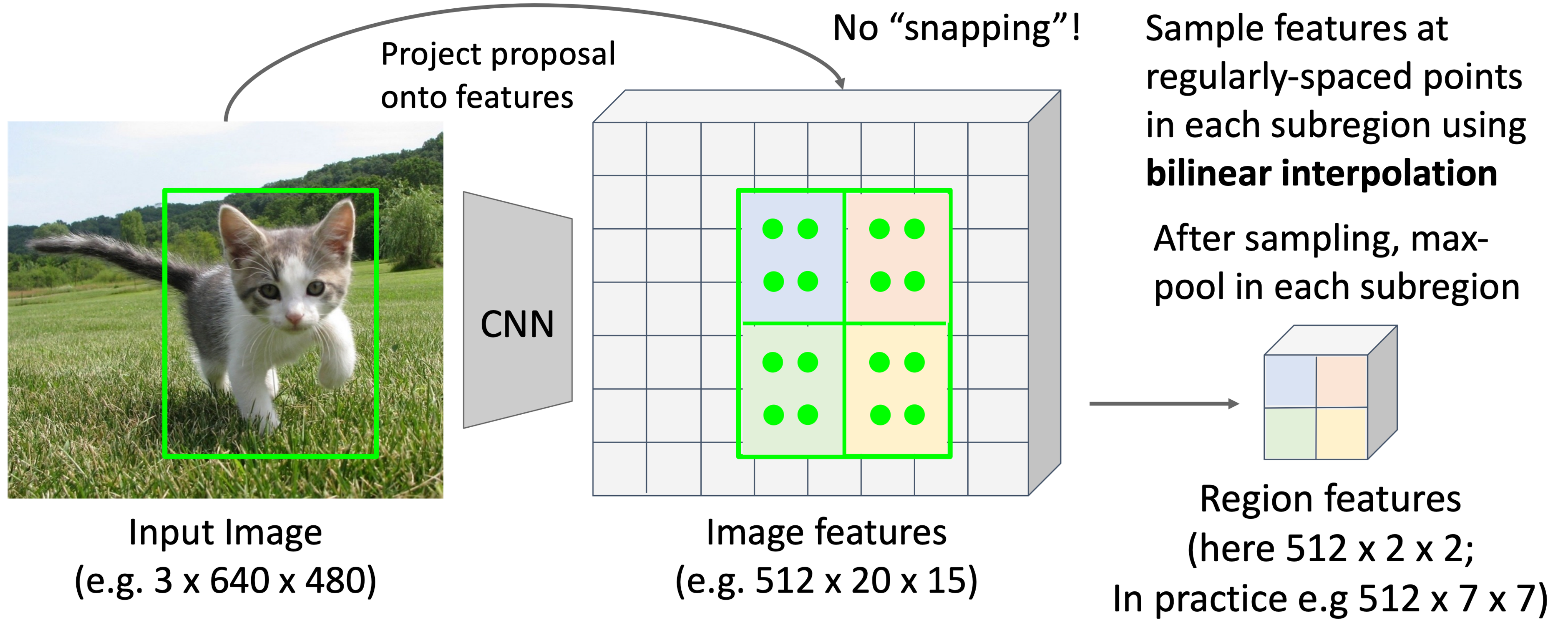
$$f_{6.5,5.8} = (f_{6,5} * 0.5 * 0.2) + (f_{7,5} * 0.5 * 0.2) + (f_{6,6} * 0.5 * 0.8) + (f_{7,6} * 0.5 * 0.8)$$

Feature f_{xy} for point (x, y) is a linear combination of features at its four neighboring grid cells:

He et al, "Mask R-CNN", ICCV 2017



Cropping Features: RoI Align



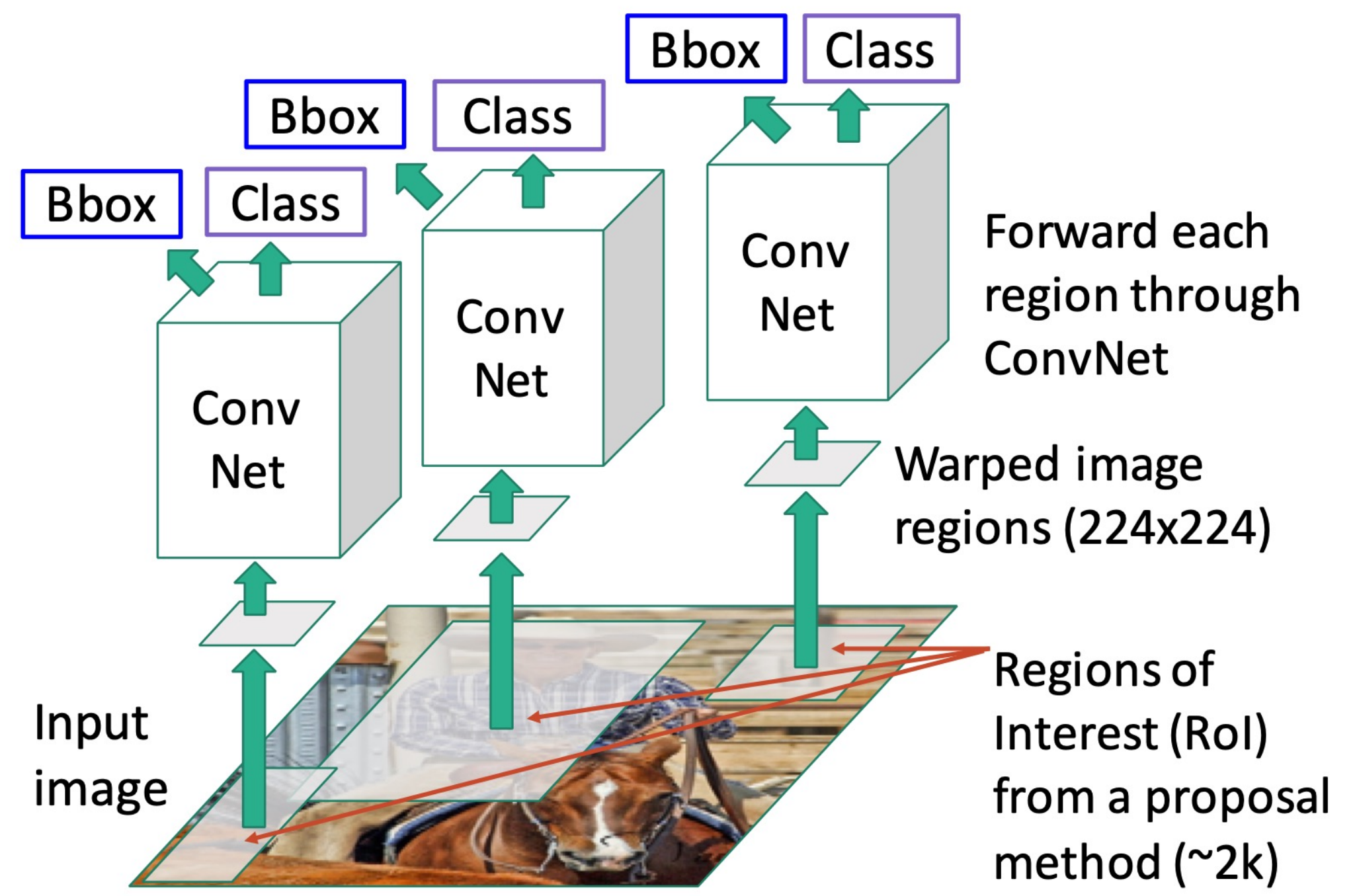
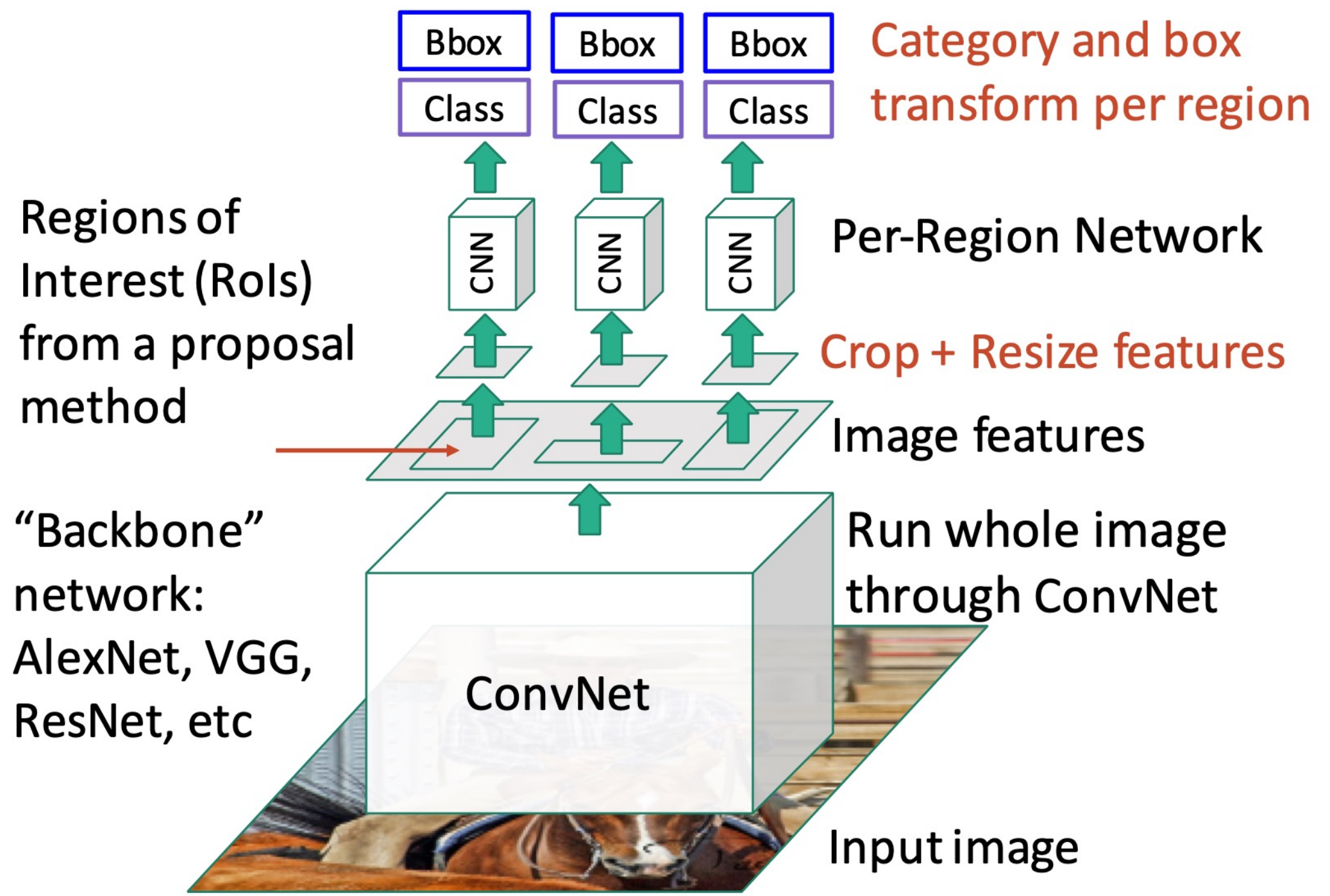
He et al, "Mask R-CNN", ICCV 2017



Fast R-CNN vs "Slow" R-CNN

Fast R-CNN: Apply differentiable cropping to shared image features

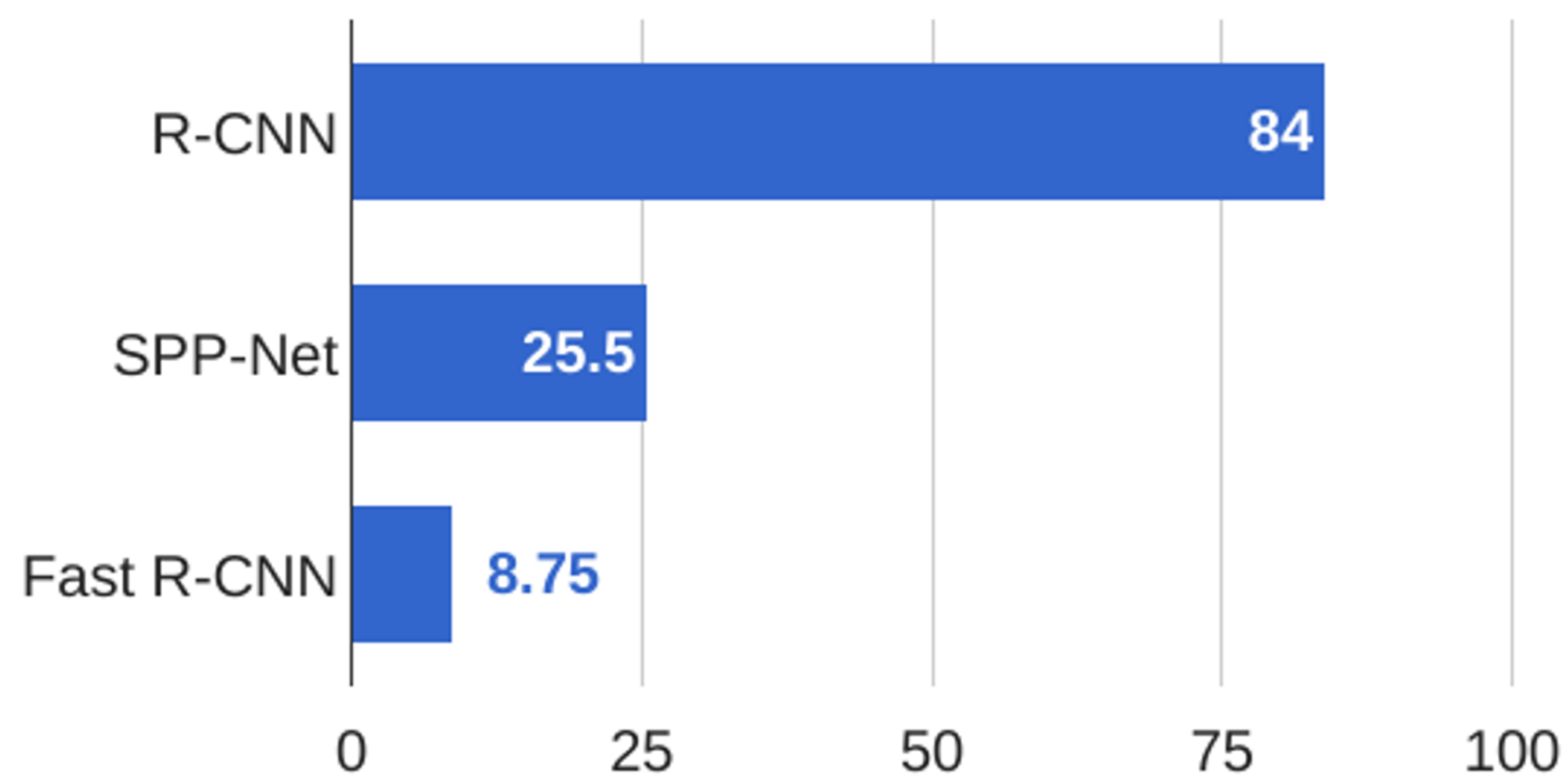
"Slow" R-CNN: Apply differentiable cropping to shared image features



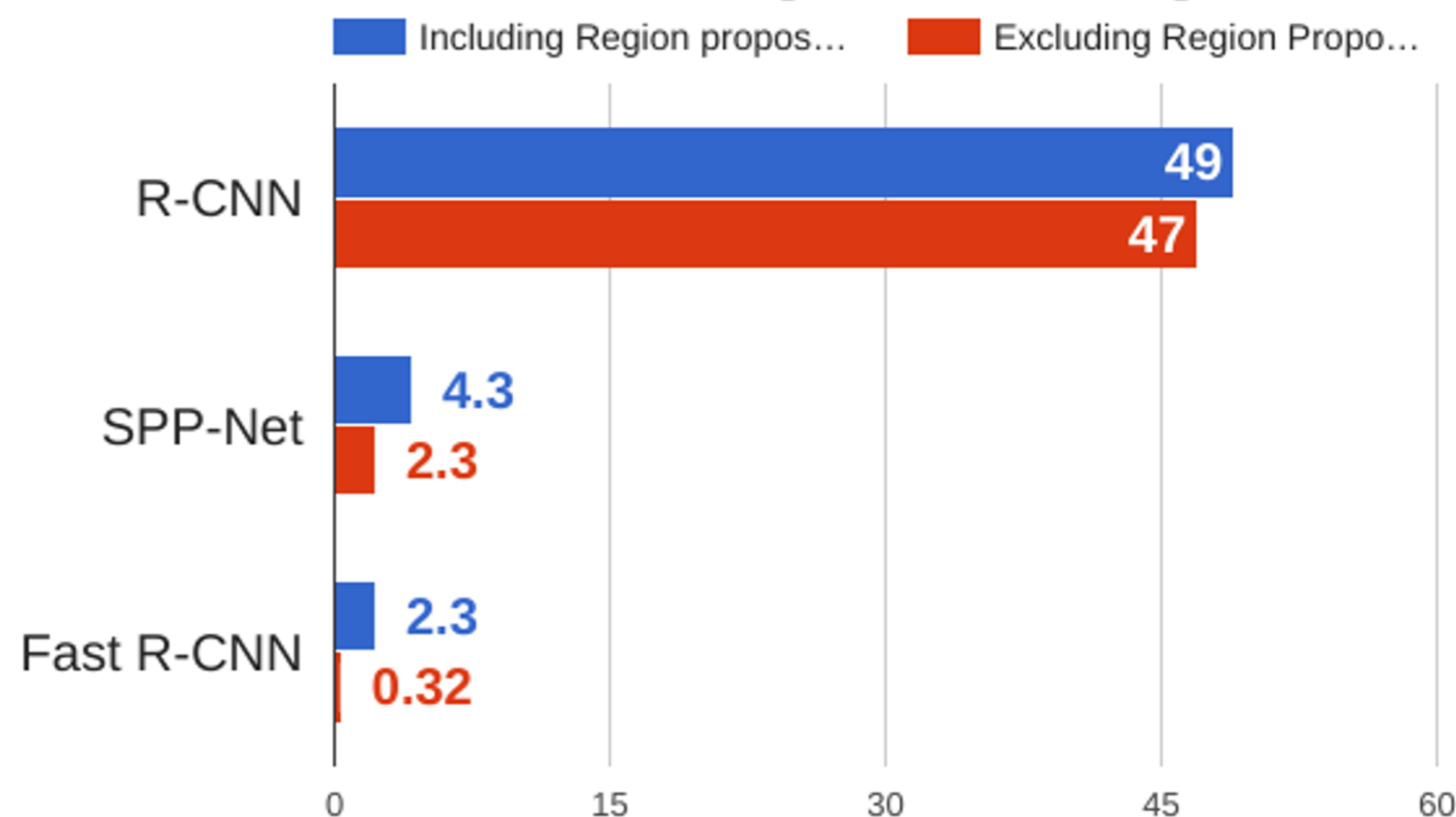


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)

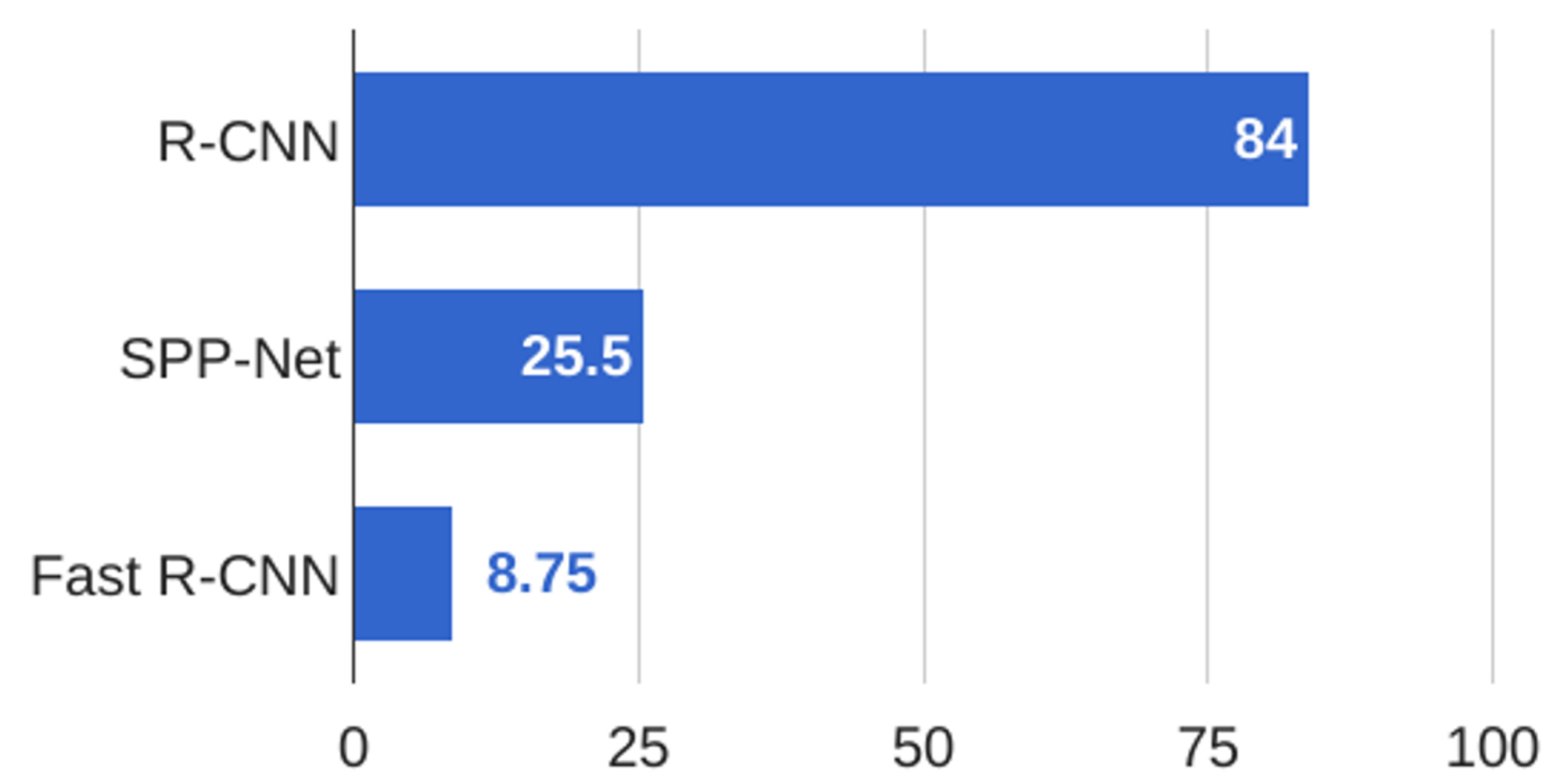


Girshick et al, “Rich feature hierarchies for accurate object detection and

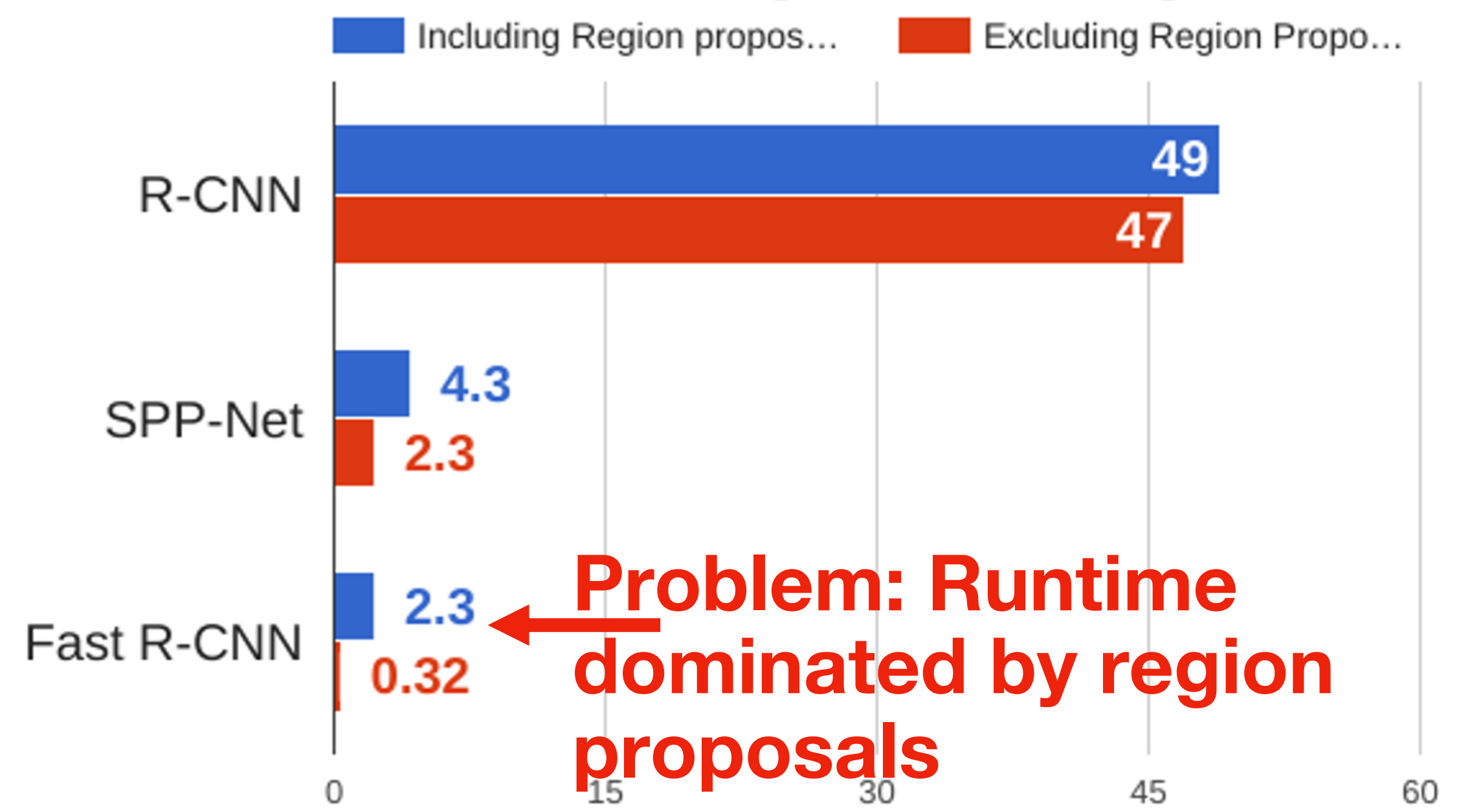


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)



Problem: Runtime dominated by region proposals

Girshick et al, “Rich feature hierarchies for accurate object detection and

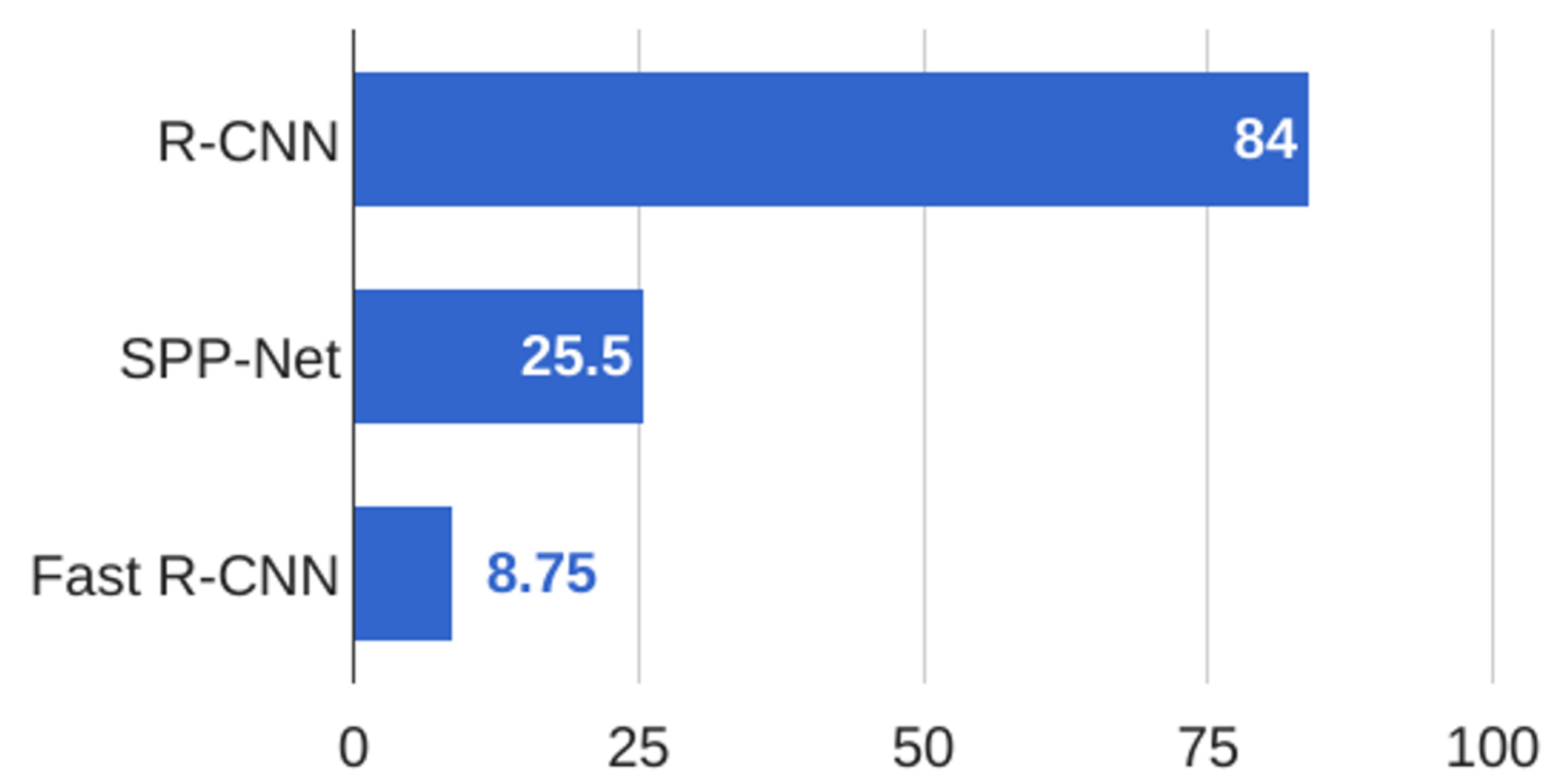


He et al, “Spatial pyramid pooling in deep convolutional networks for vi

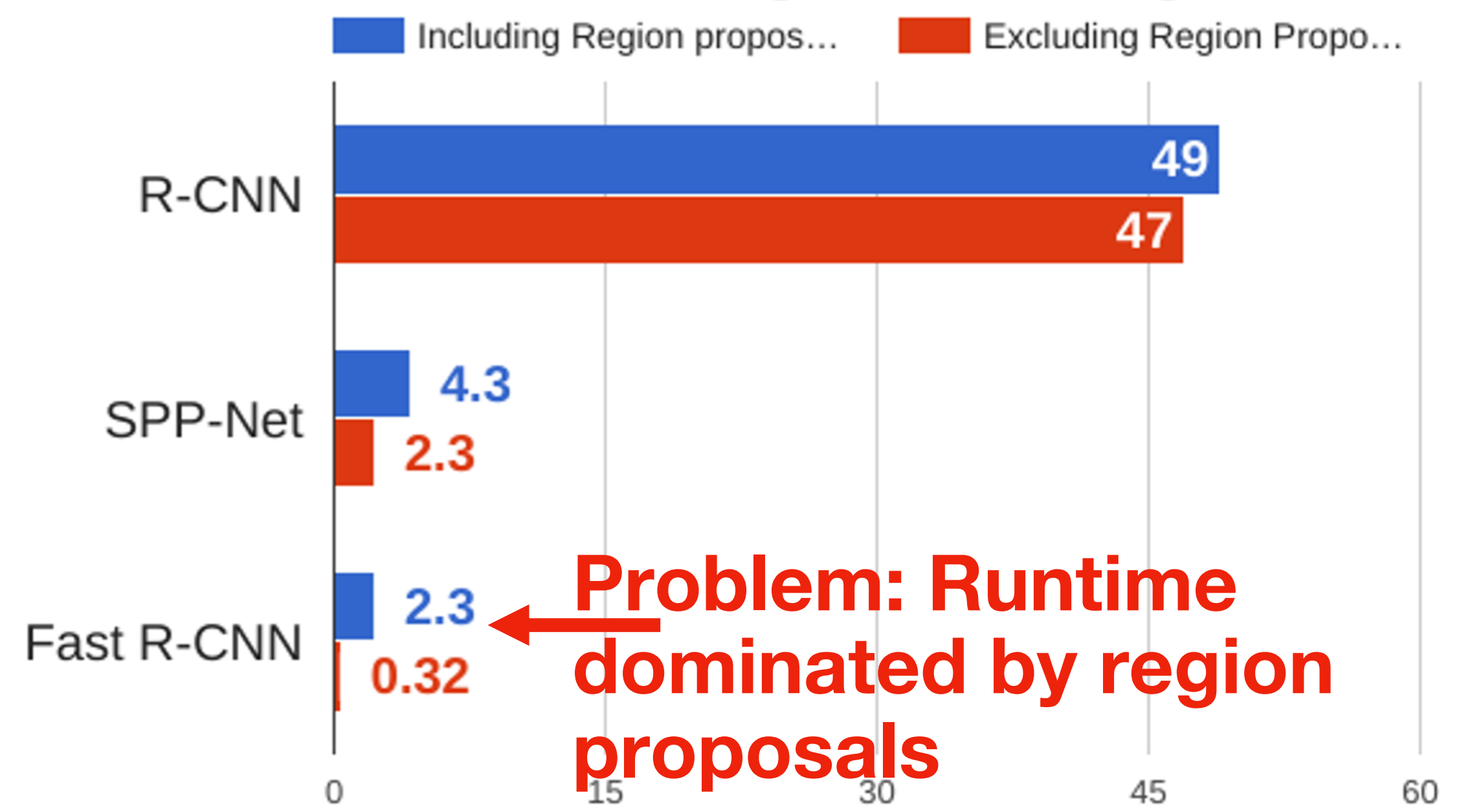


Fast R-CNN vs “Slow” R-CNN

Training time (Hours)



Test time (seconds)



Problem: Runtime dominated by region proposals

Recall: Region proposals computed by heuristic “Selective search” algorithm on CPU – let’s learn them with a CNN

Girshick et al, “Rich feature hierarchies for accurate object detection and semantic segmentation”, CVPR 2014.

He et al, “Spatial pyramid pooling in deep convolutional networks for visual recognition”, ECCV 2014
 Girshick, “Fast R-CNN”, ICCV 2015

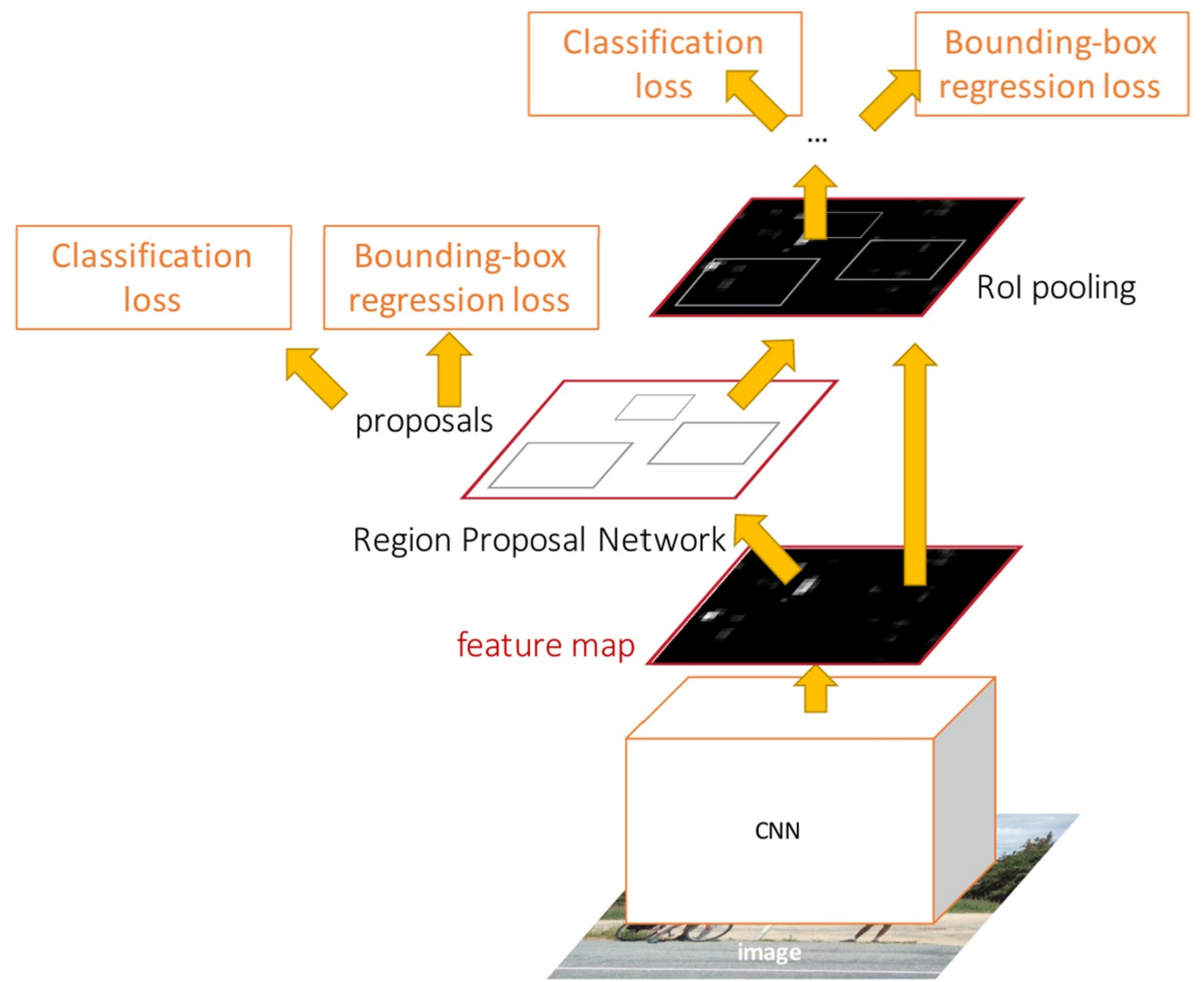




Faster R-CNN: Learnable Region Proposals

Insert **Region Proposal Network (RPN)** to predict proposals from features

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one





Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)

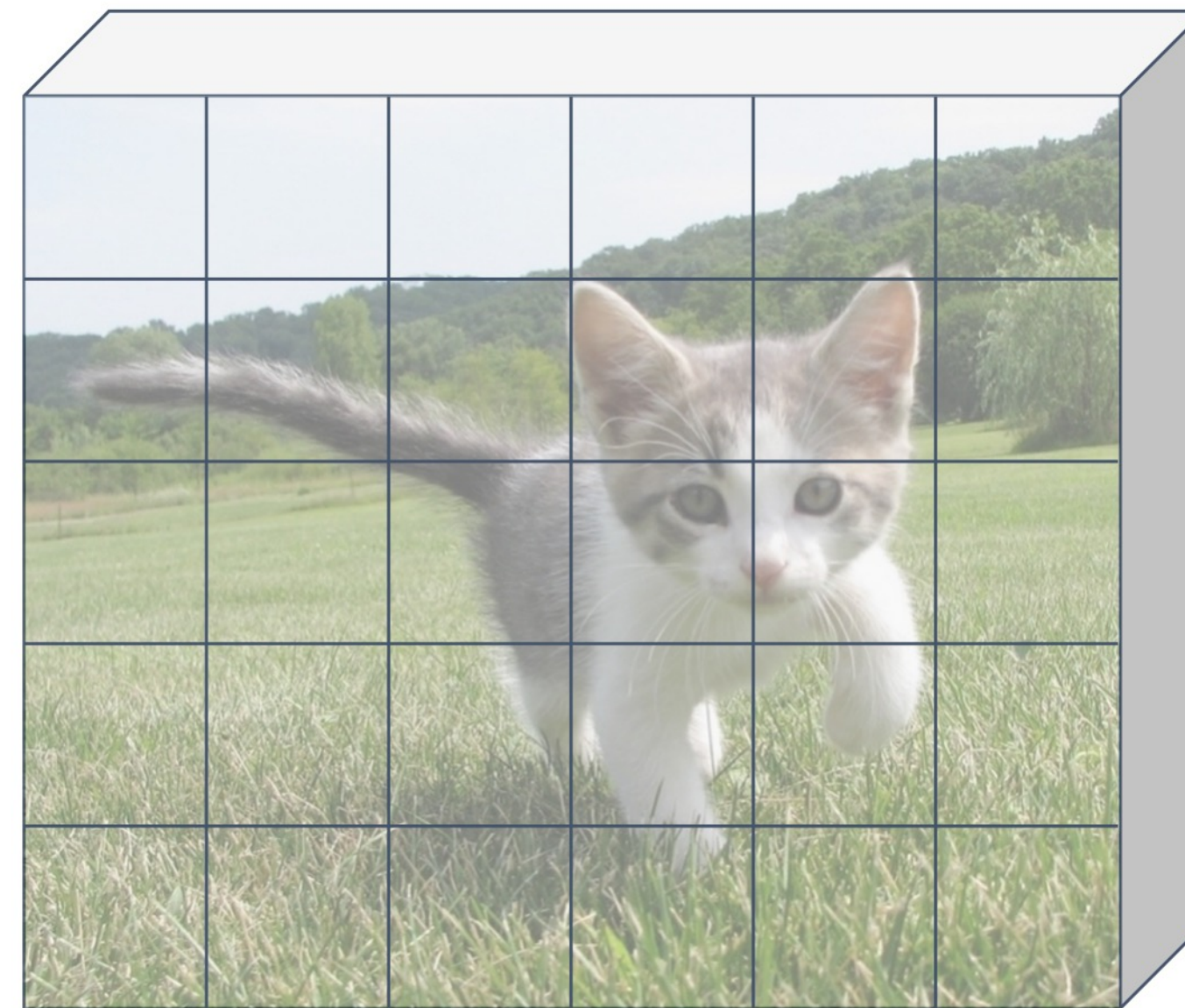


Image features
(e.g. 512 x 5 x 6)



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

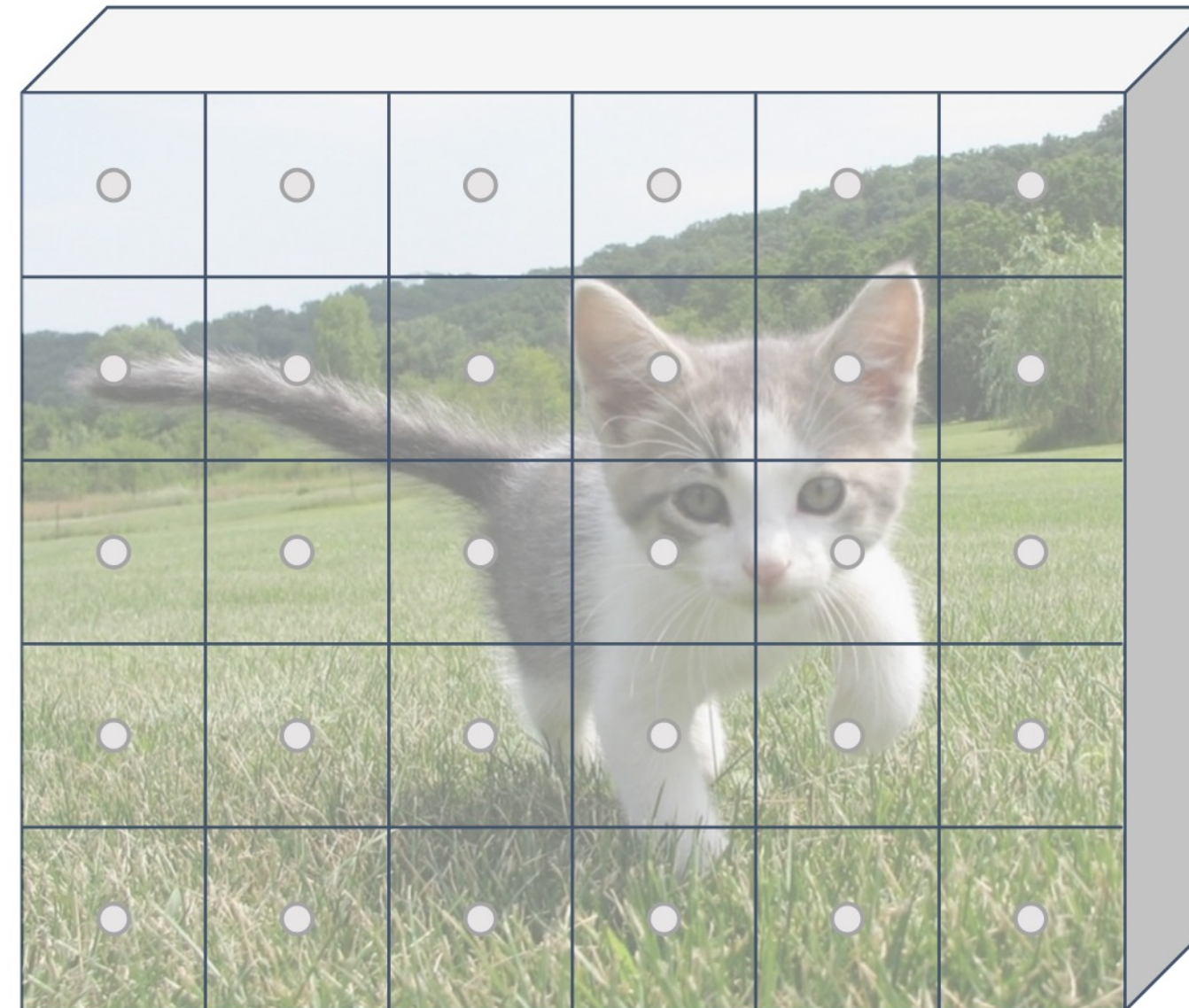
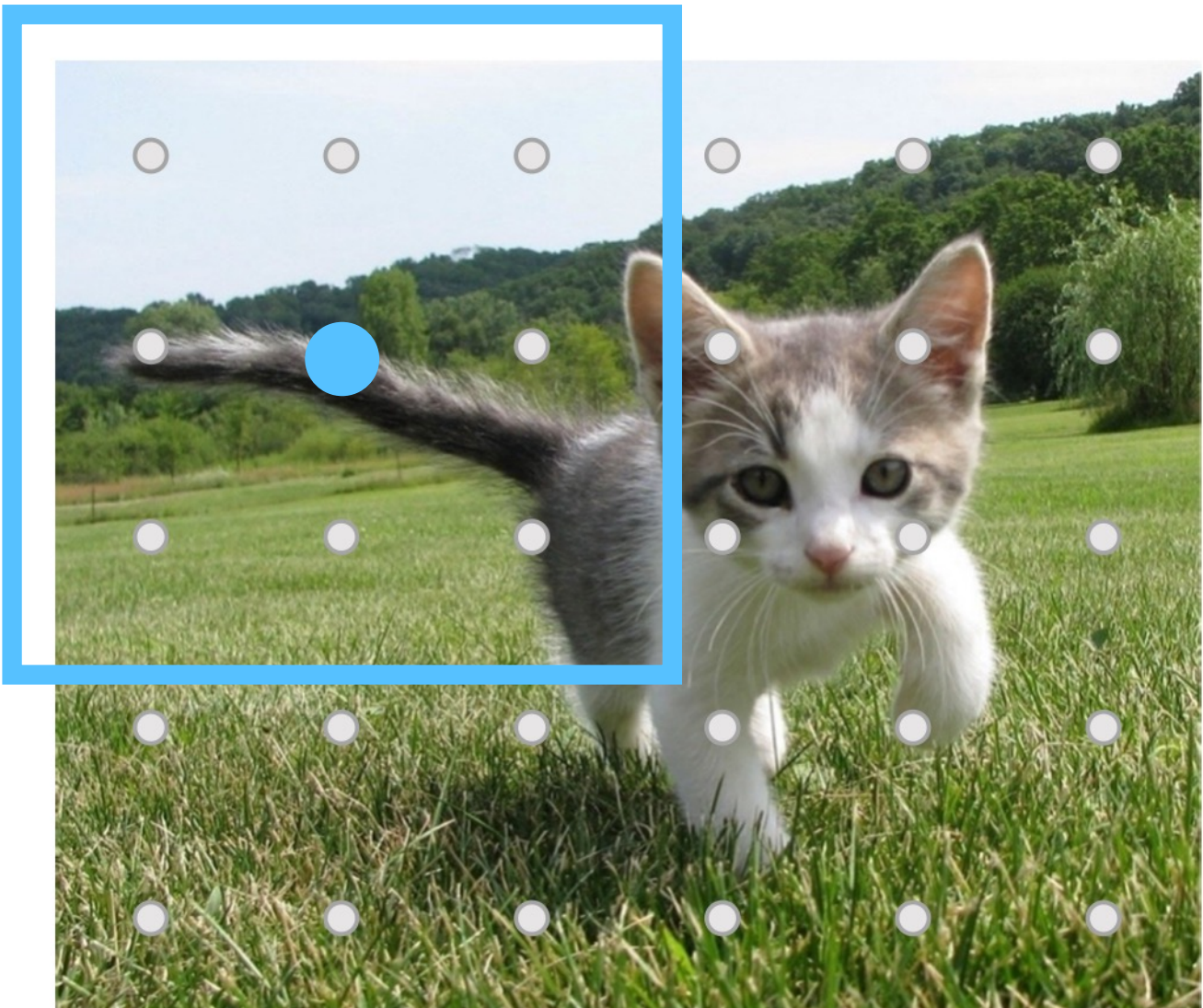


Image features
(e.g. 512 x 5 x 6)



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

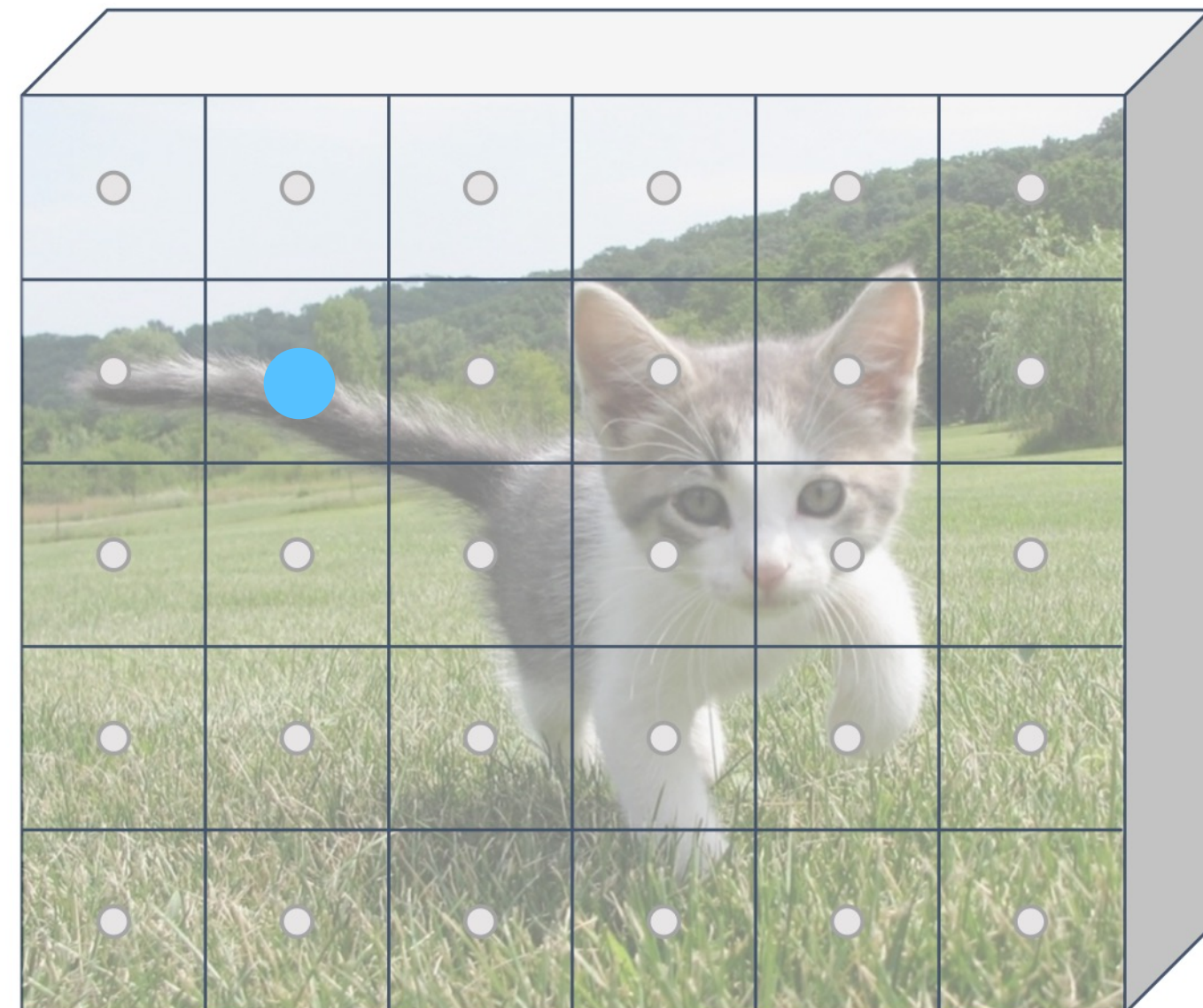


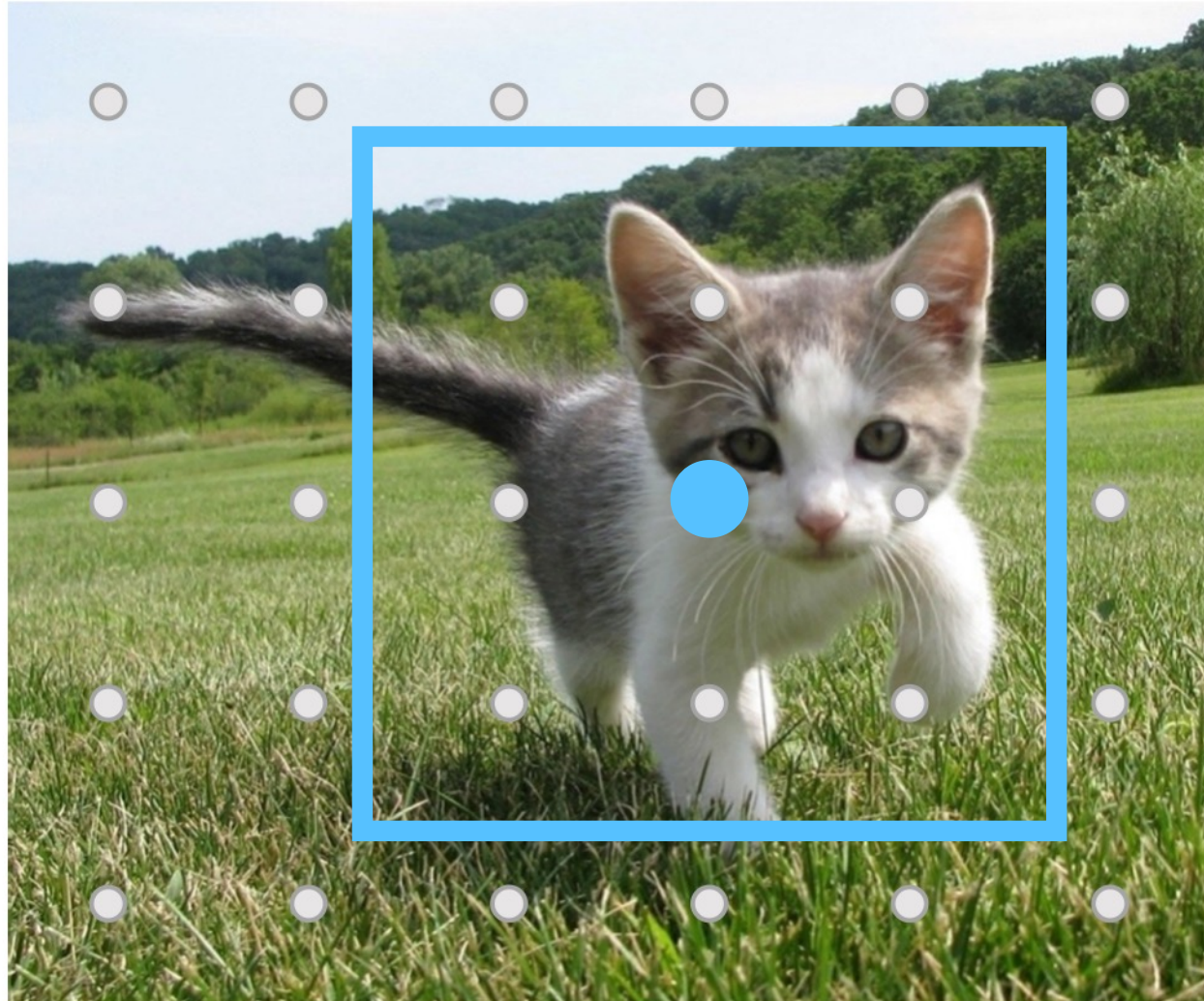
Image features
(e.g. 512 x 5 x 6)

Imagine an **anchor box** of fixed size at each point in the feature map



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

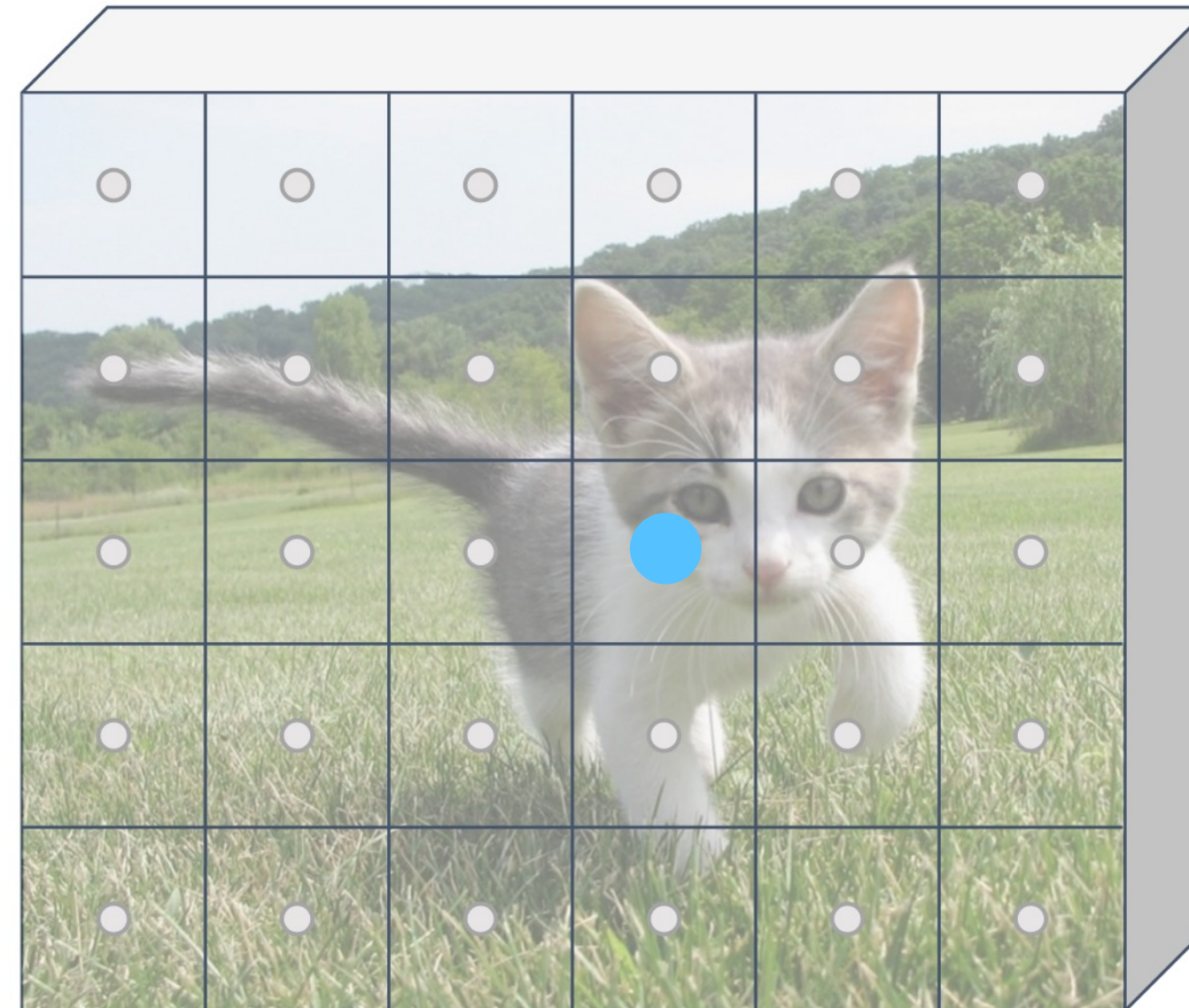


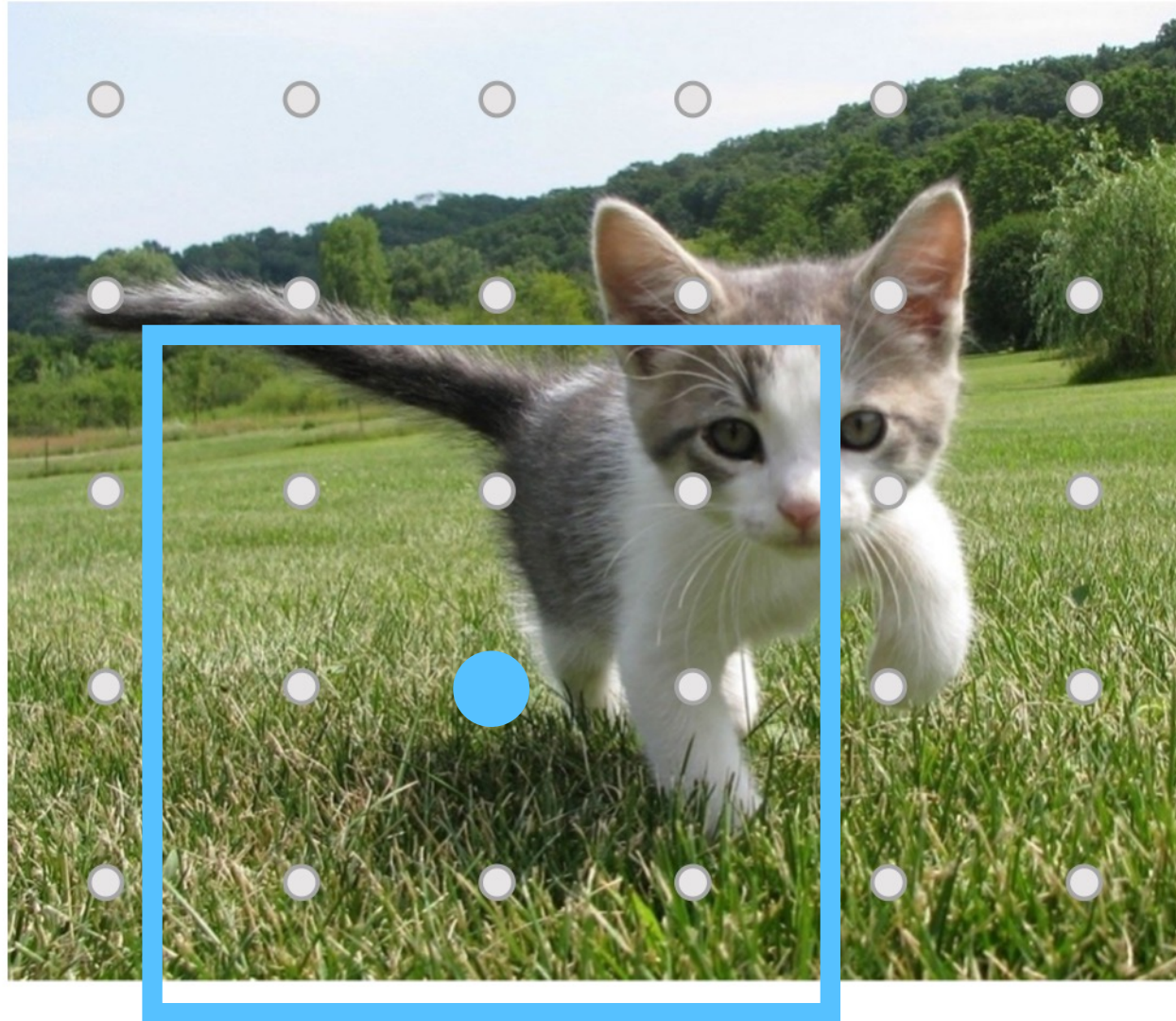
Image features
(e.g. 512 x 5 x 6)

Imagine an **anchor box** of fixed size at each point in the feature map



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

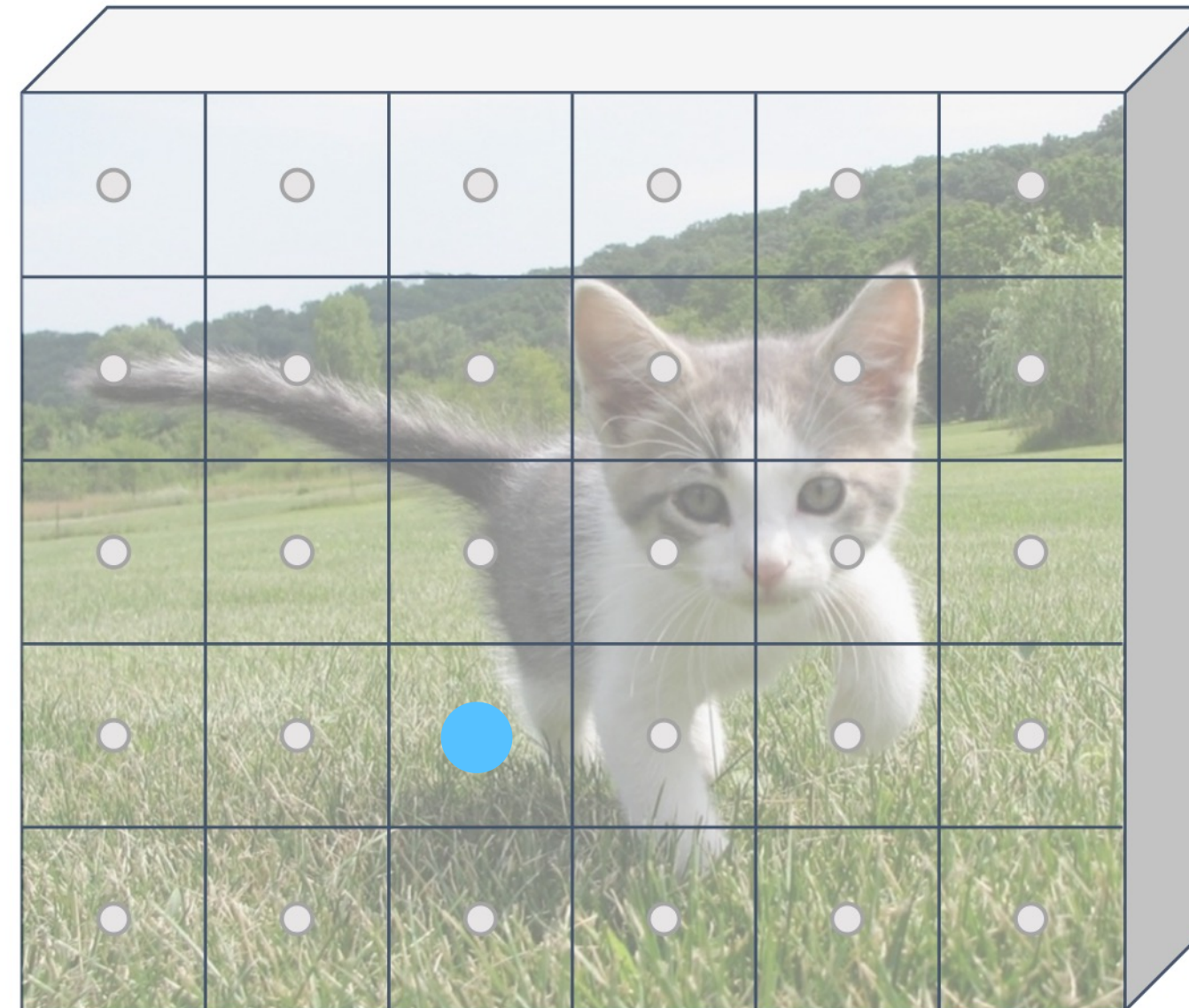


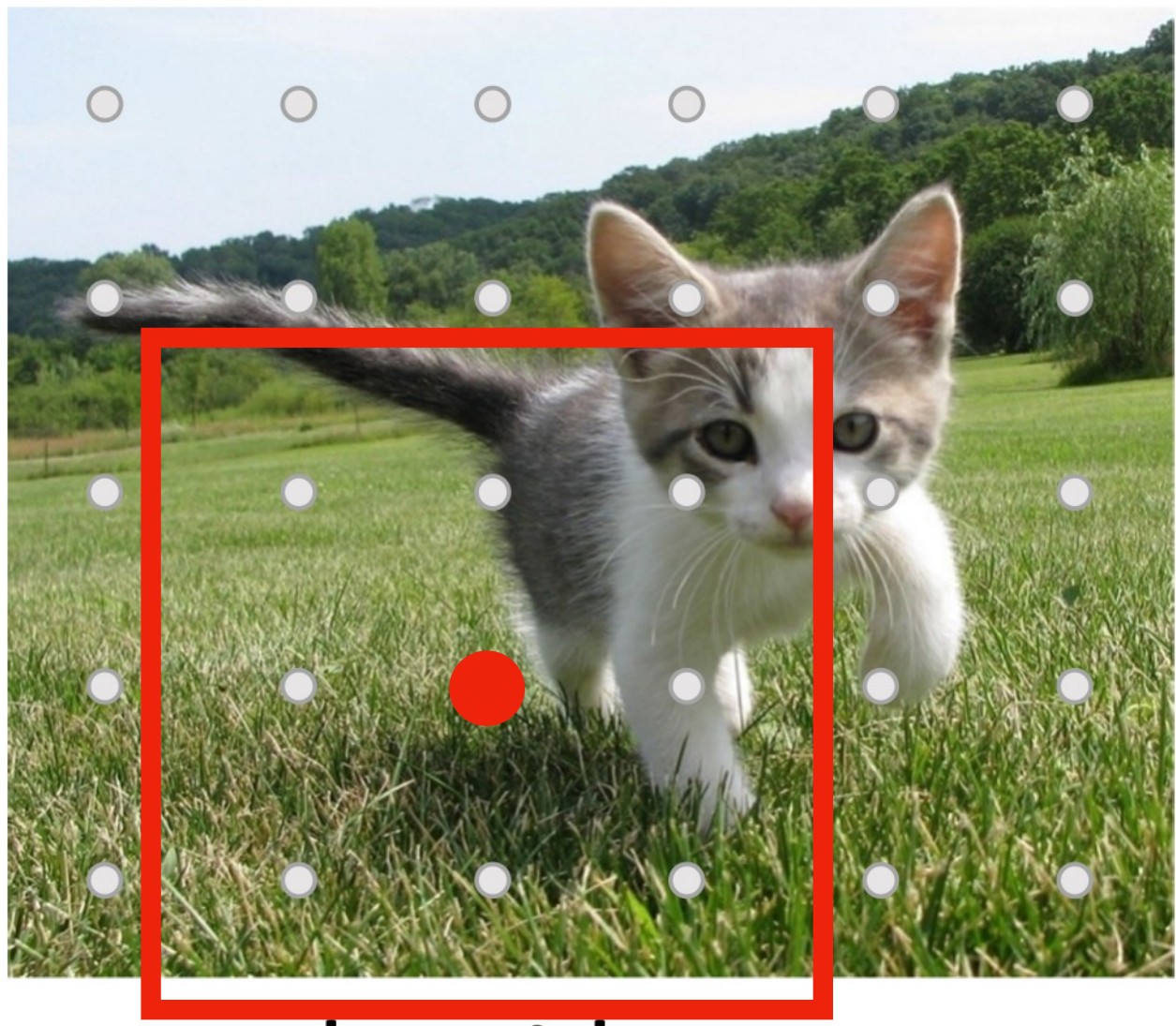
Image features
(e.g. 512 x 5 x 6)

Imagine an **anchor box** of fixed size at each point in the feature map



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

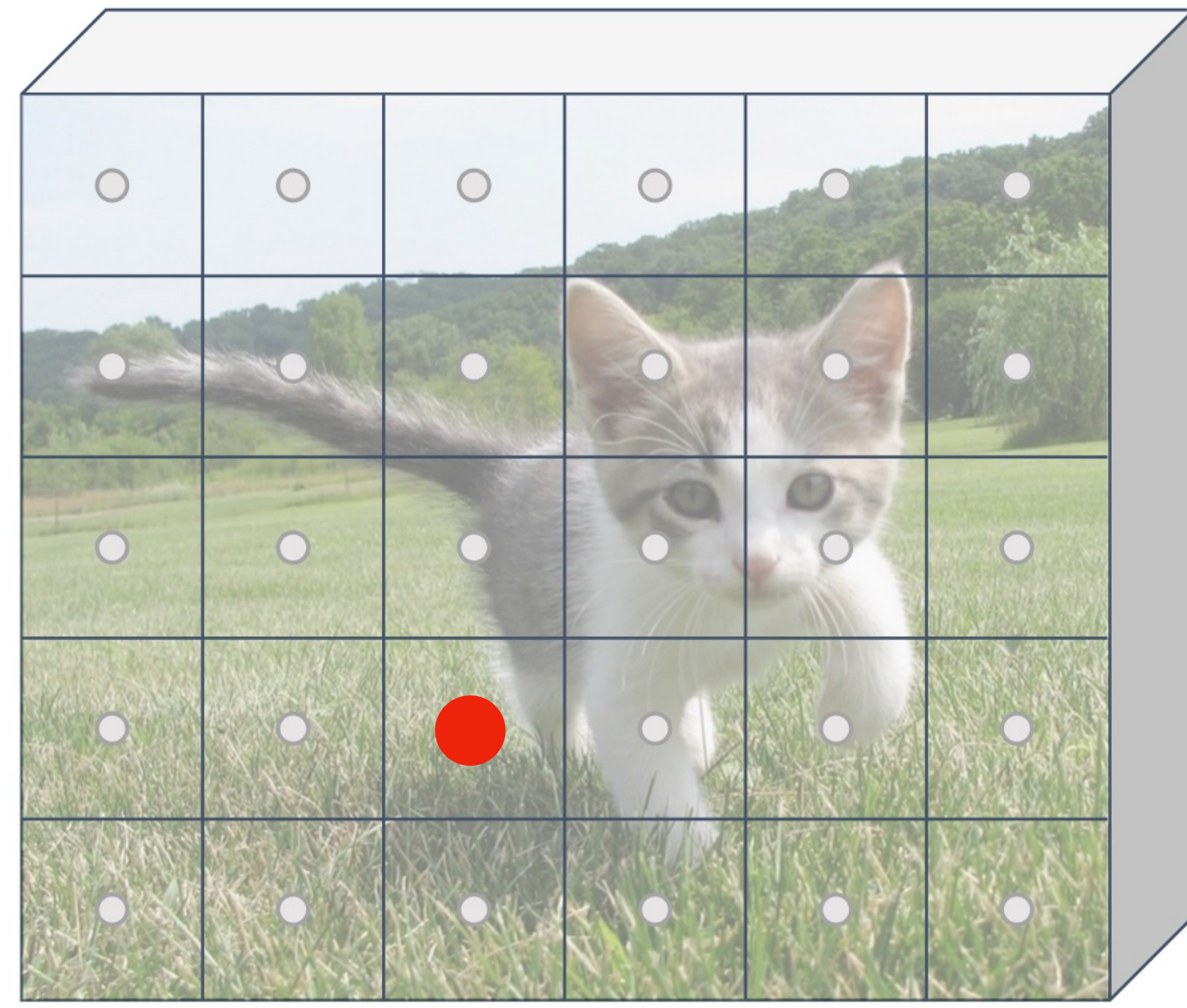


Image features
(e.g. 512 x 5 x 6)

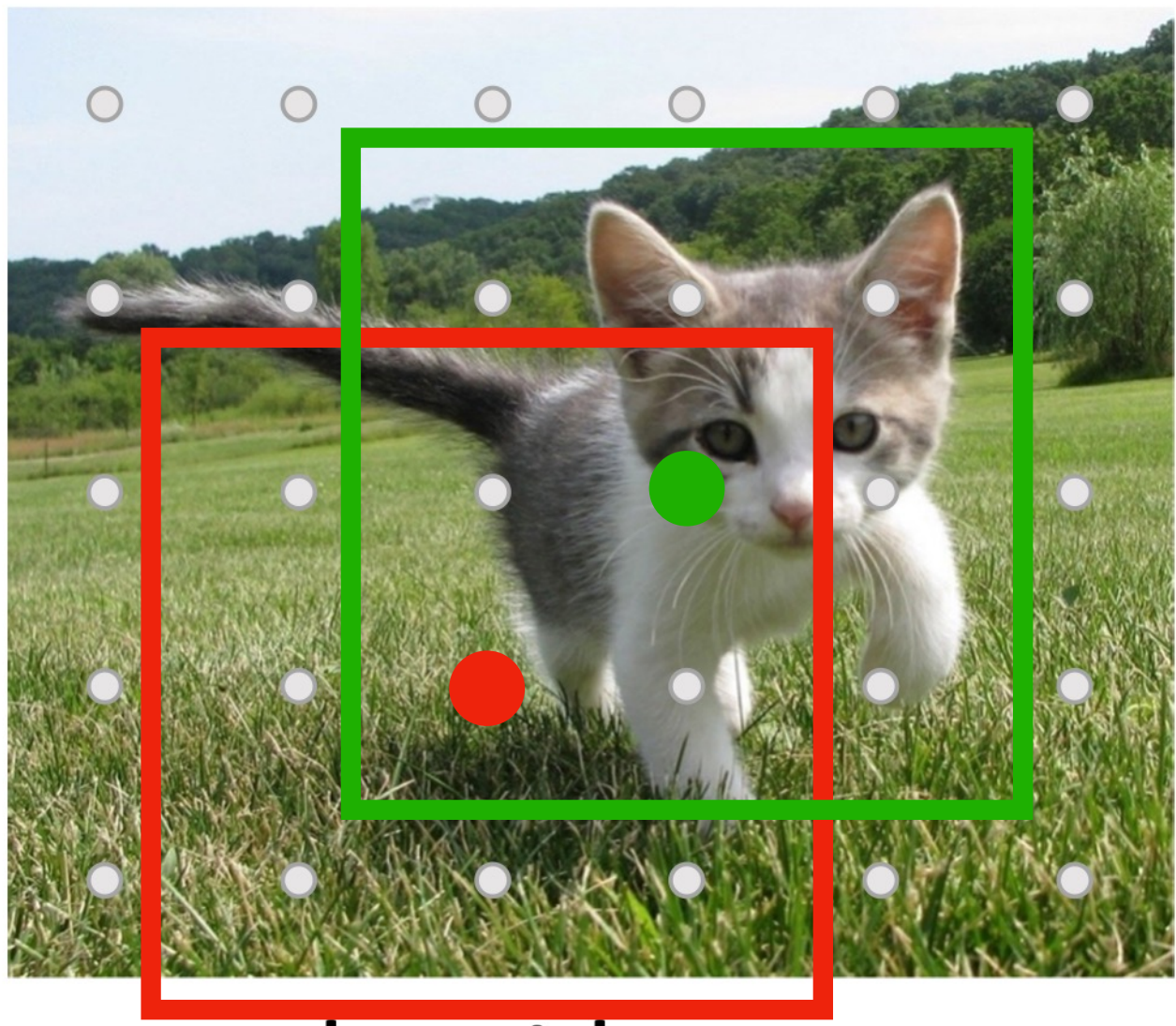
Imagine an **anchor box** of fixed size at each point in the feature map

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

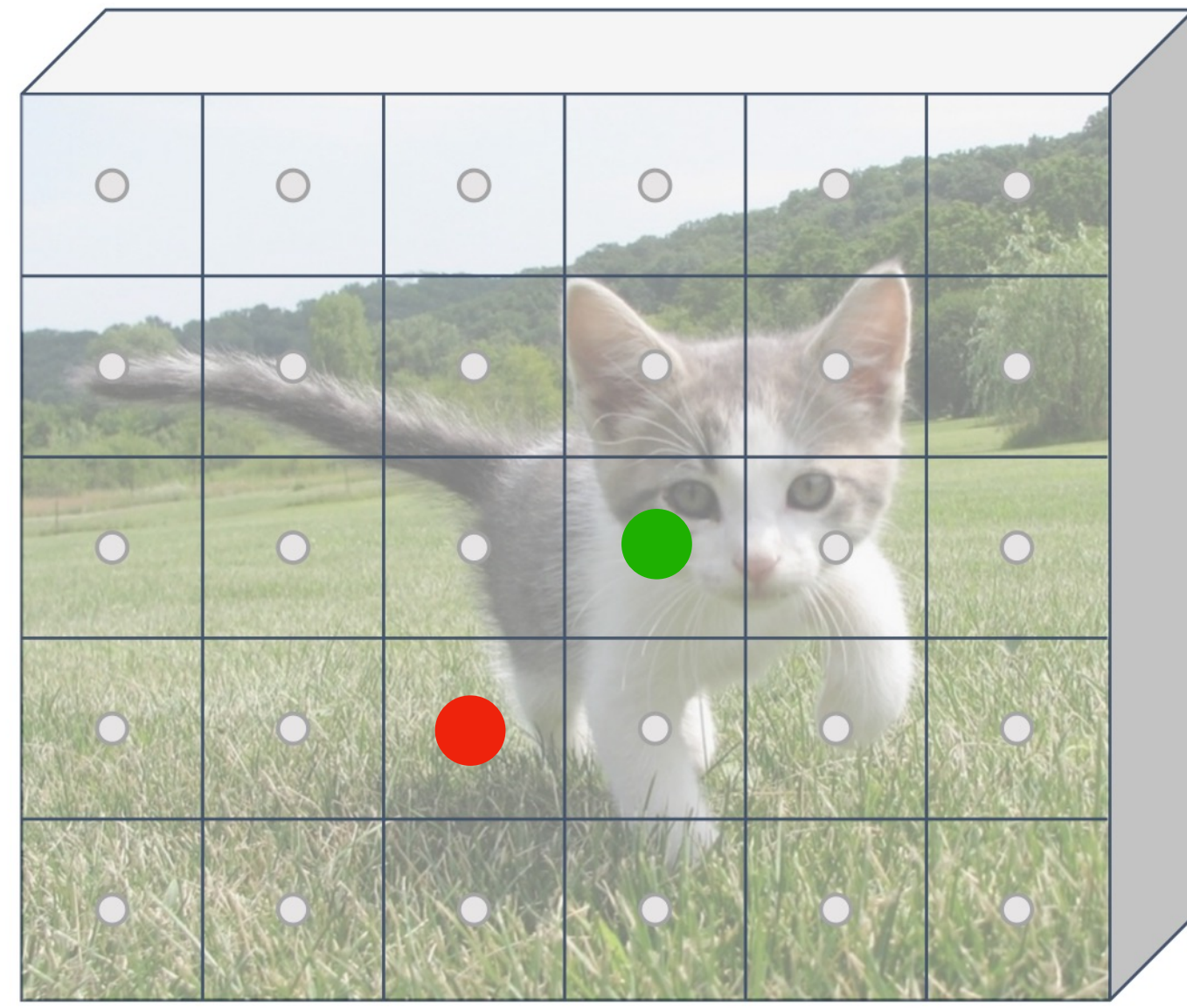


Image features
(e.g. 512 x 5 x 6)

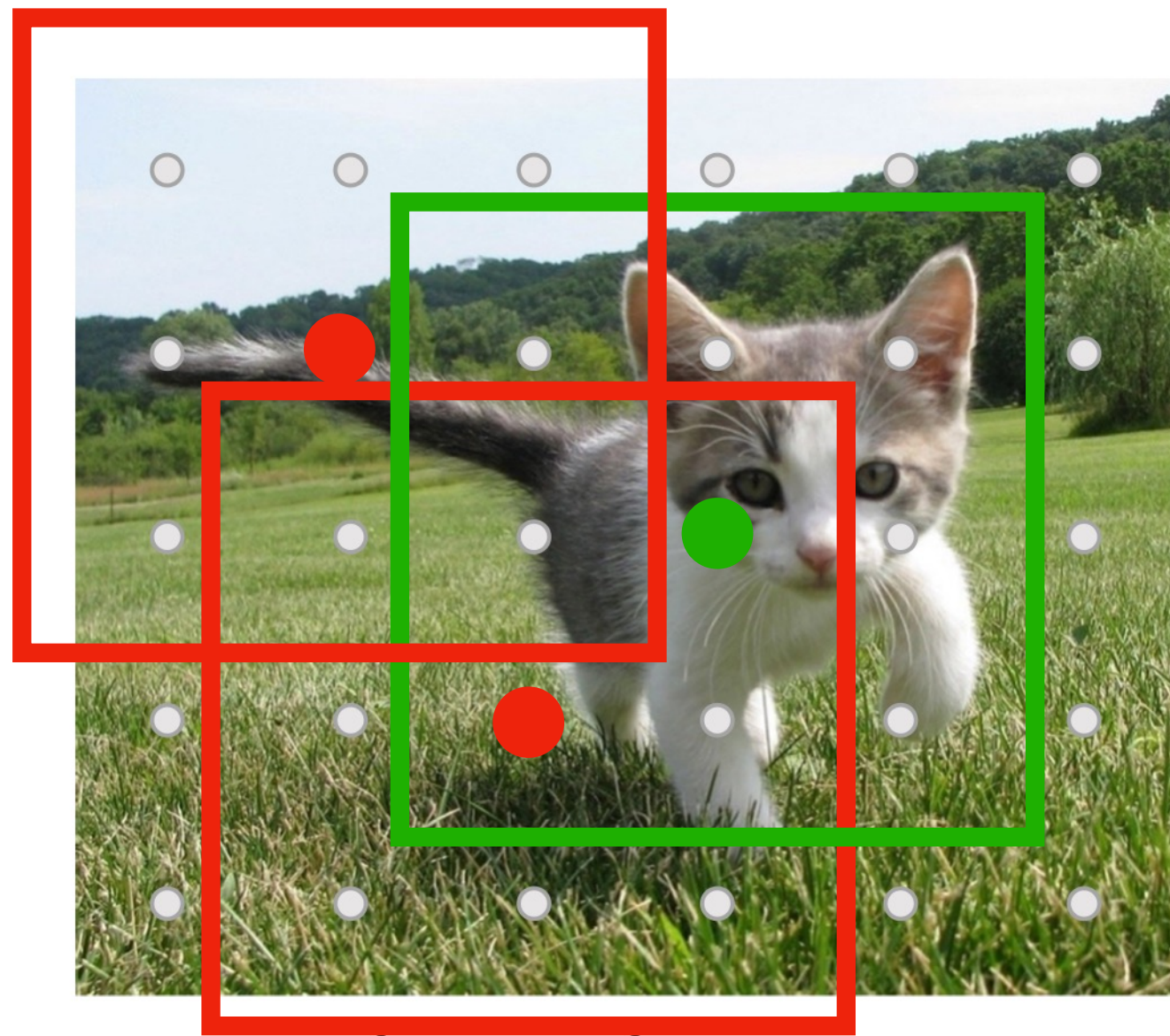
Imagine an **anchor box** of fixed size at each point in the feature map

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

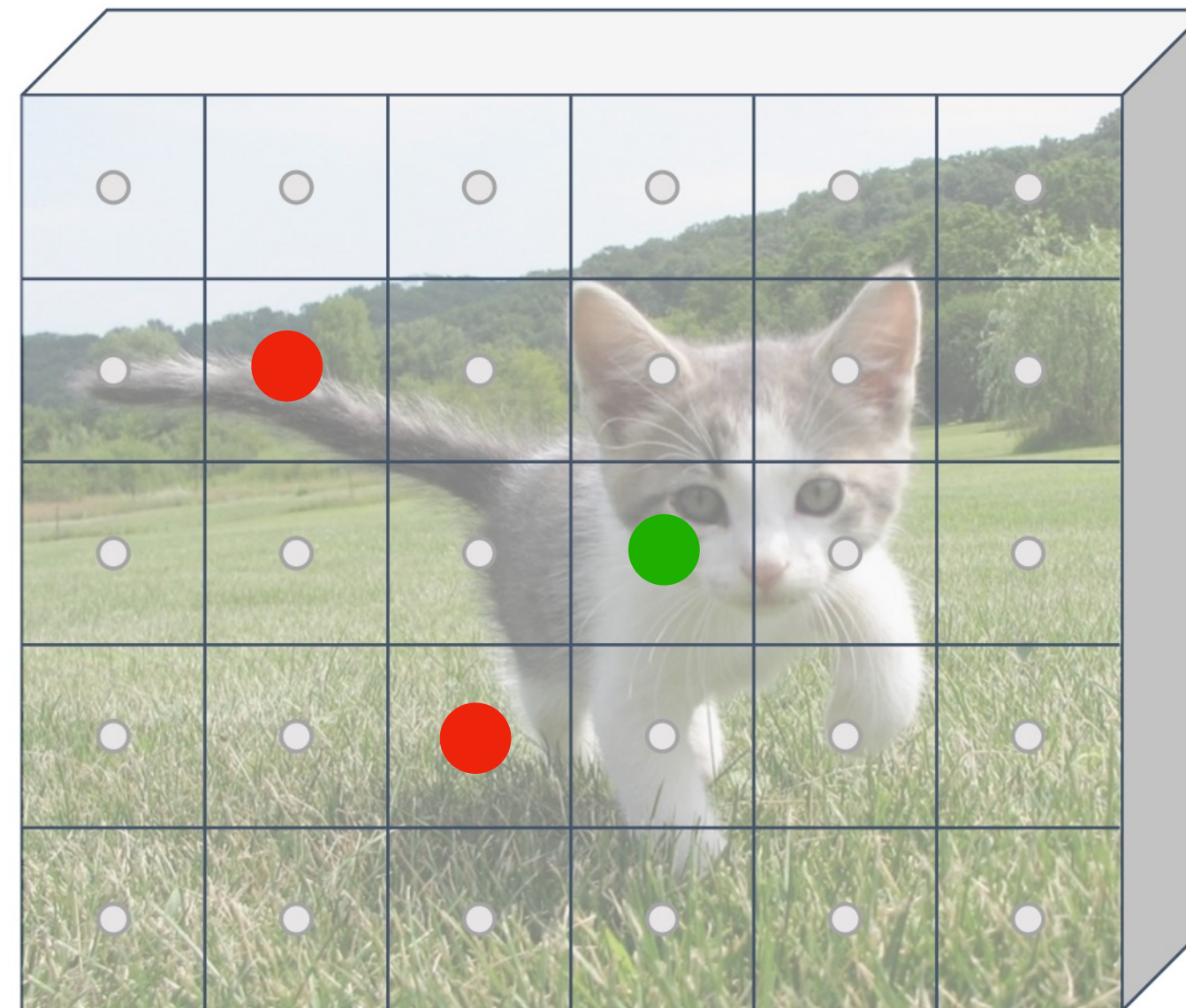


Image features
(e.g. 512 x 5 x 6)

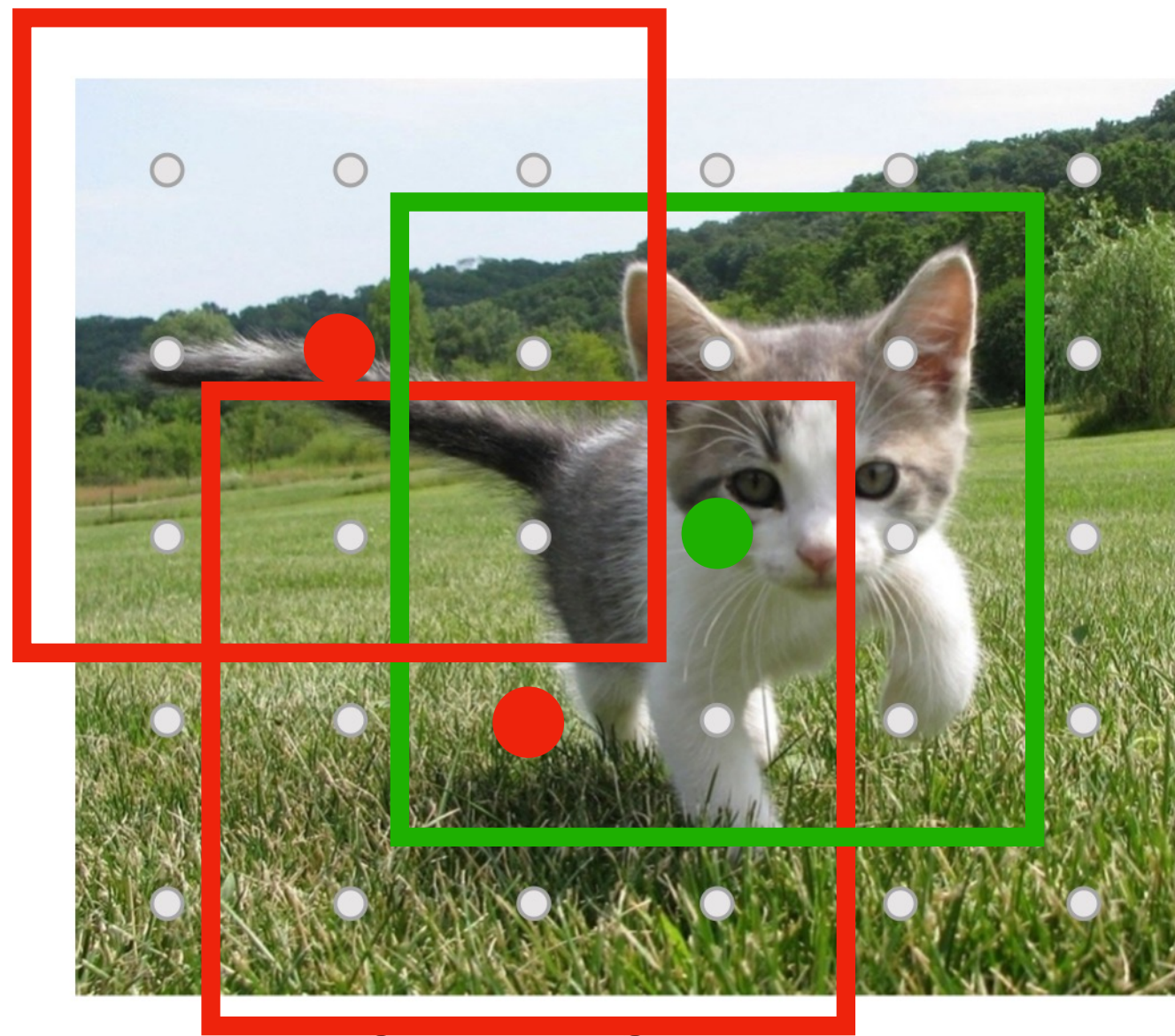
Imagine an **anchor box** of fixed size at each point in the feature map

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

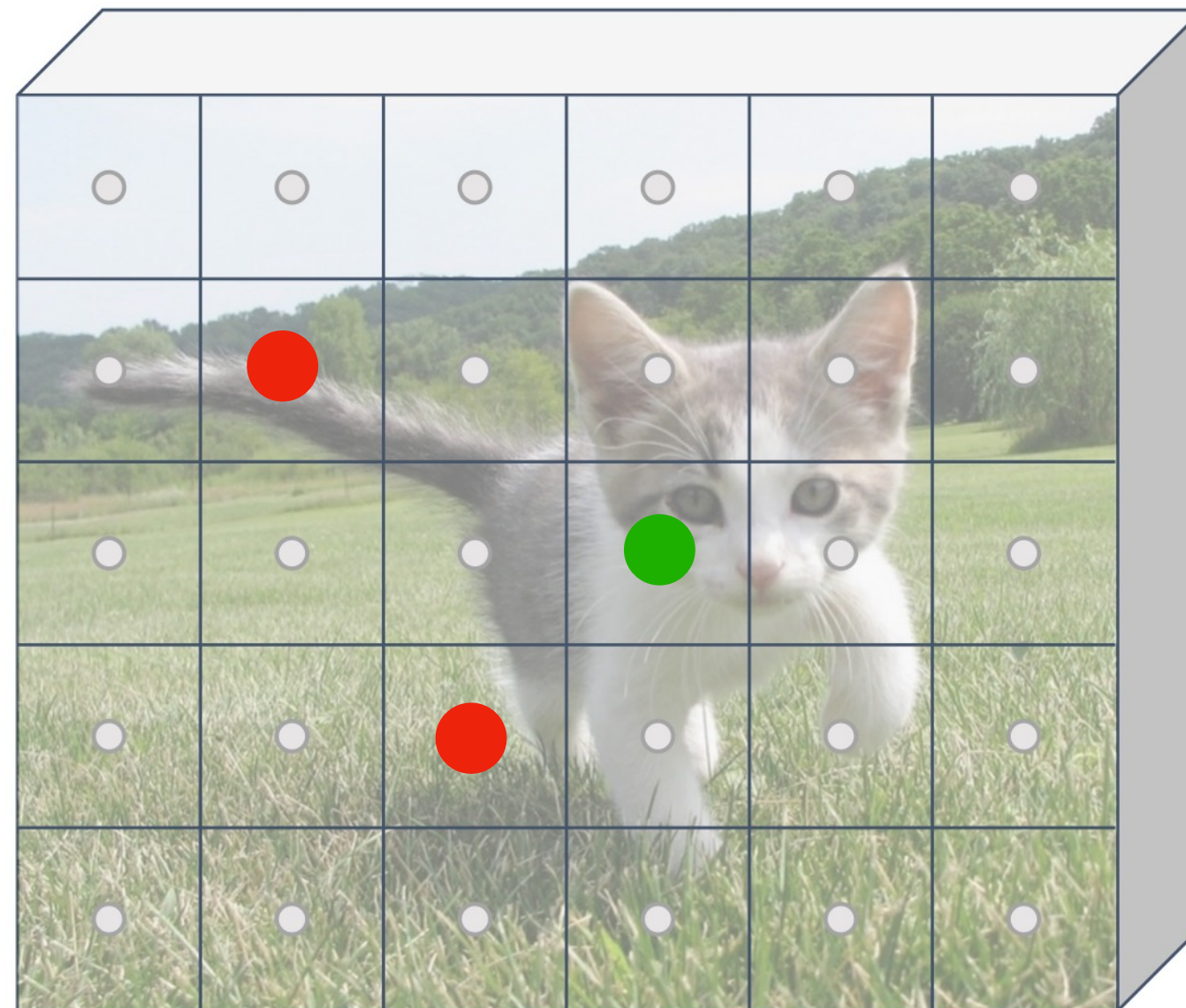


Image features
(e.g. 512 x 5 x 6)

Predict object vs not object scores for all anchors with a conv layer (512 input filters, 2 output filters)



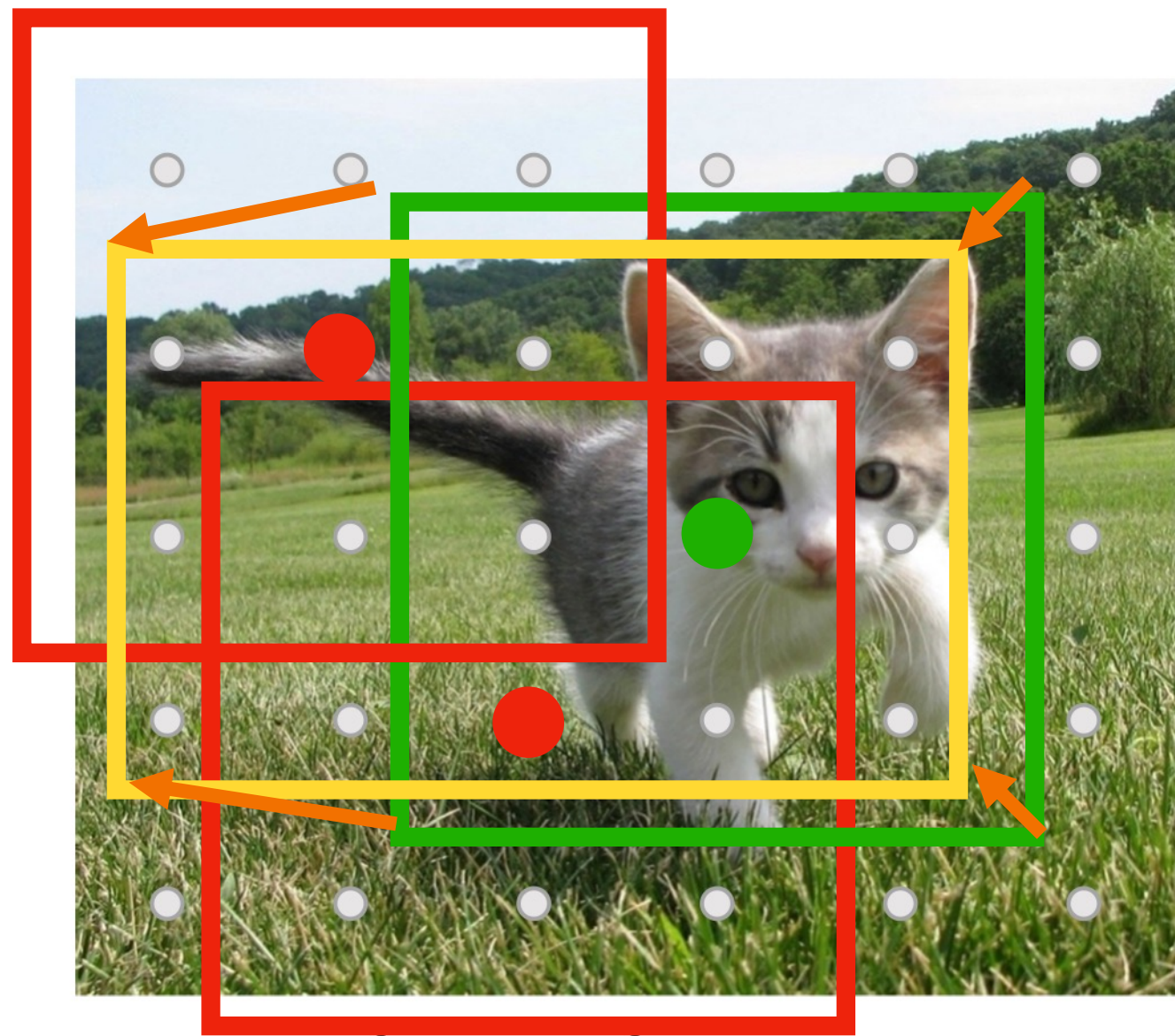
Anchor is object?
2 x 5 x 6

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

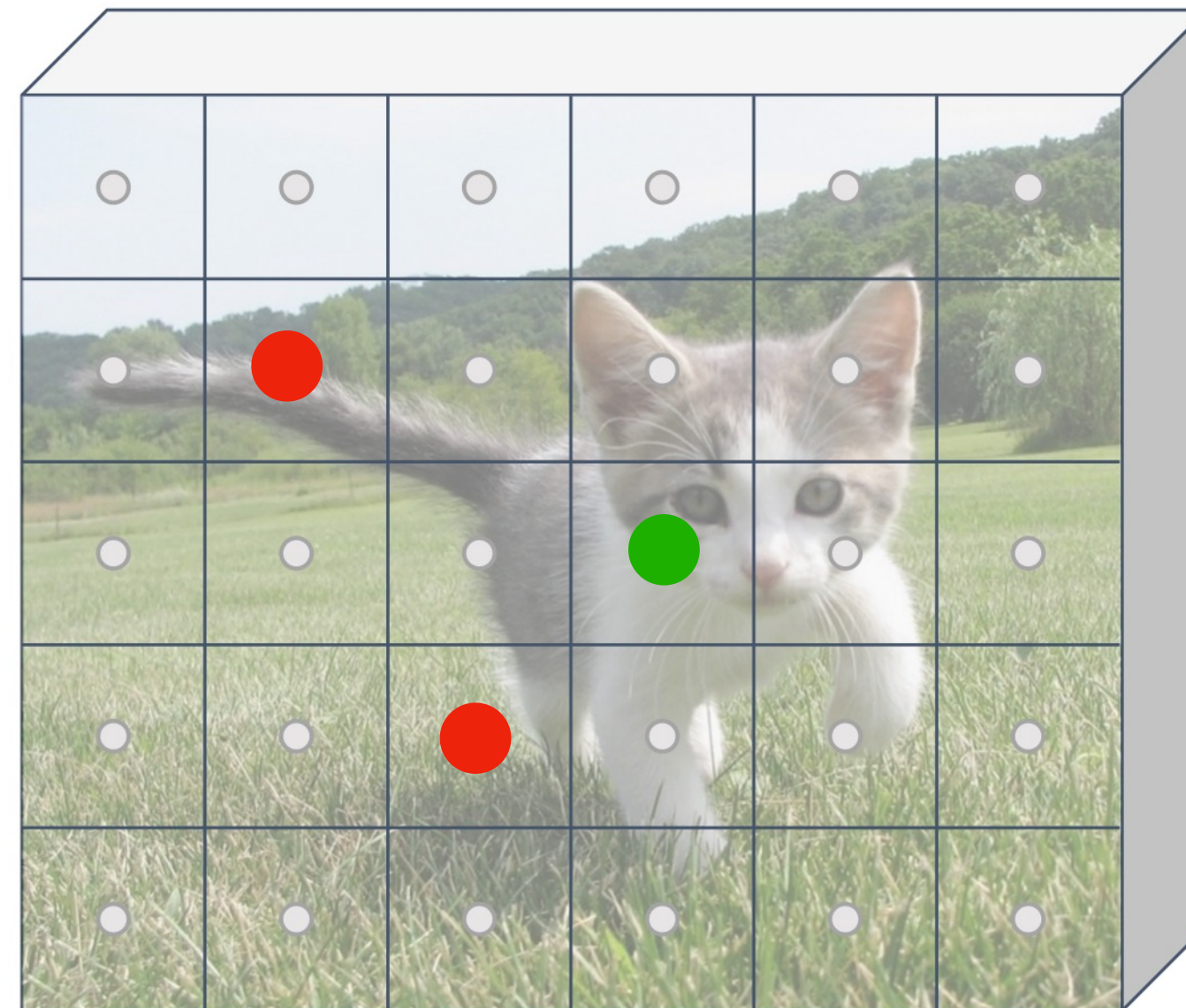


Image features
(e.g. 512 x 5 x 6)

For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)



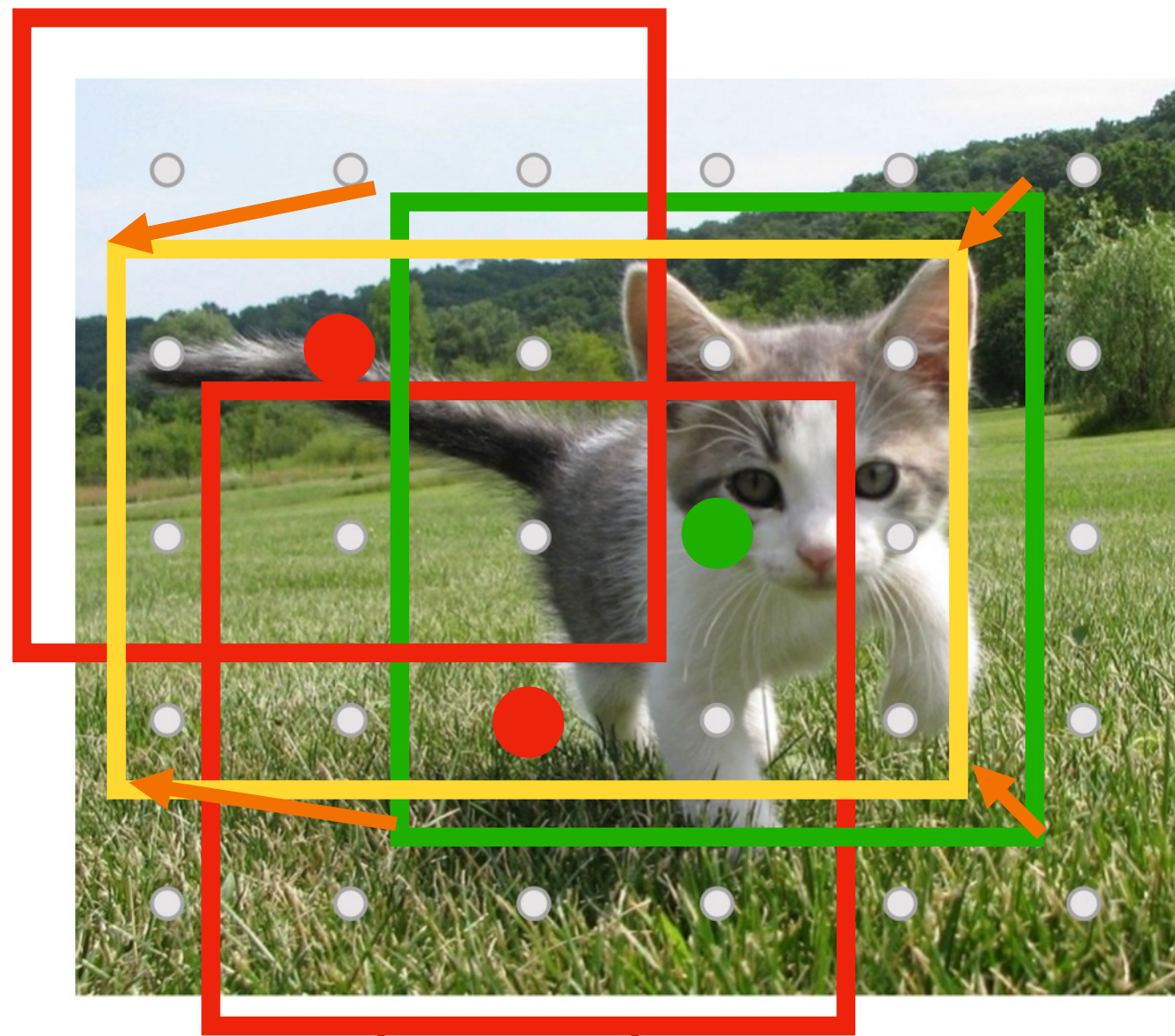
Anchor is object?
2 x 5 x 6

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

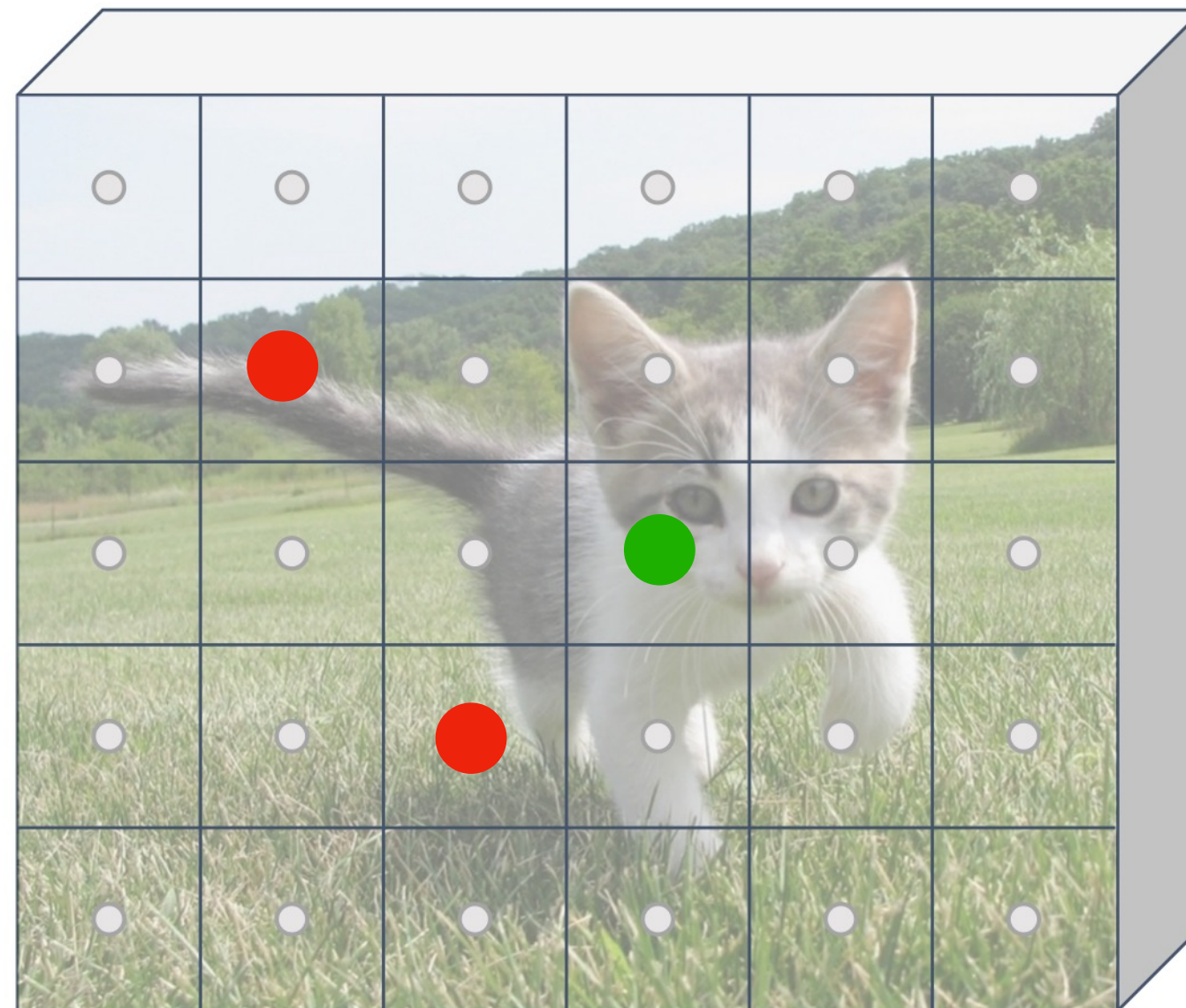


Image features
(e.g. 512 x 5 x 6)

For **positive anchors**, also predict a **transform** that converting the anchor to the **GT box** (like R-CNN)



Anchor is object?
2 x 5 x 6



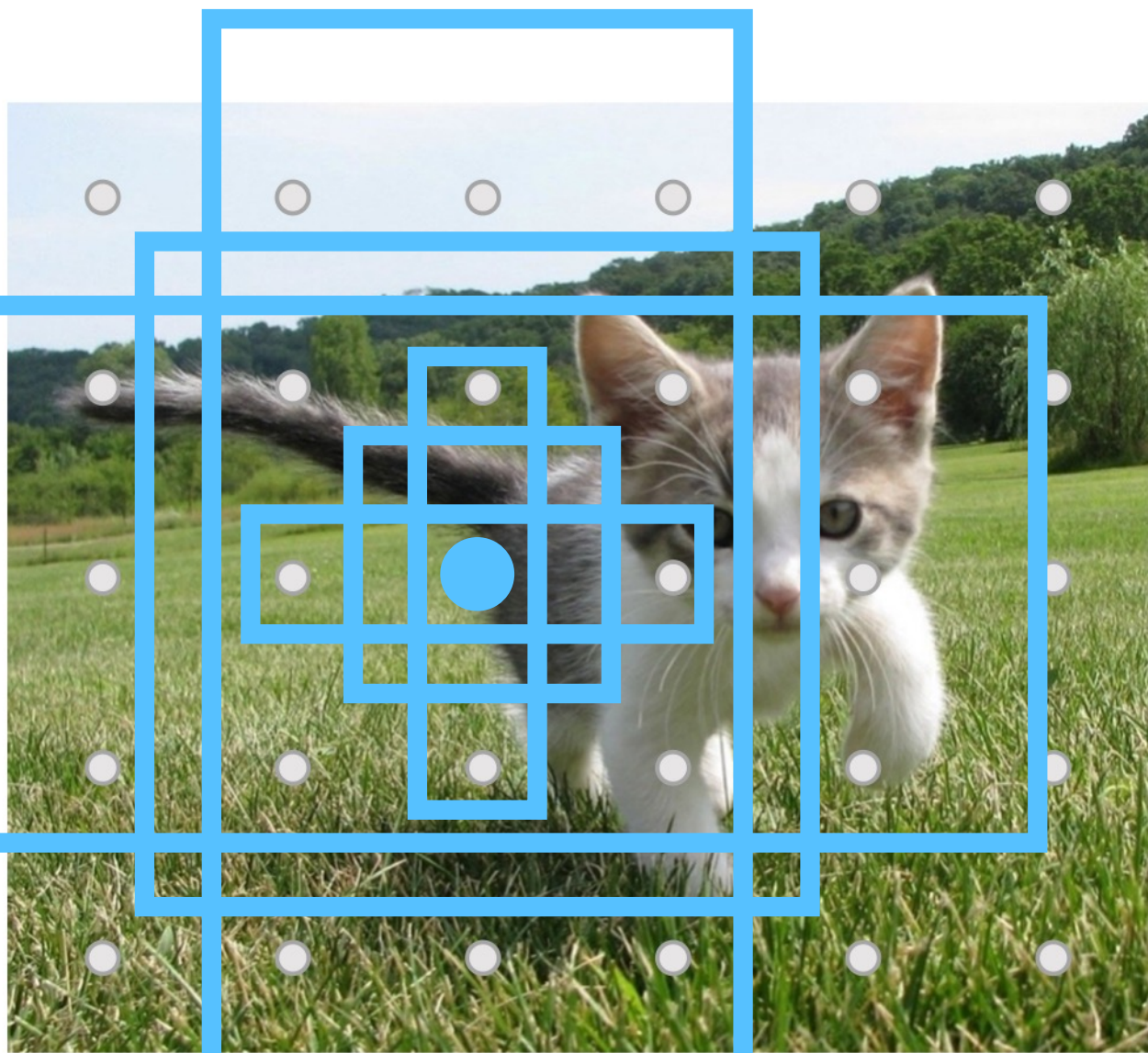
Anchor transforms
4 x 5 x 6

Classify each anchor as **positive (object)** or **negative (no object)**



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

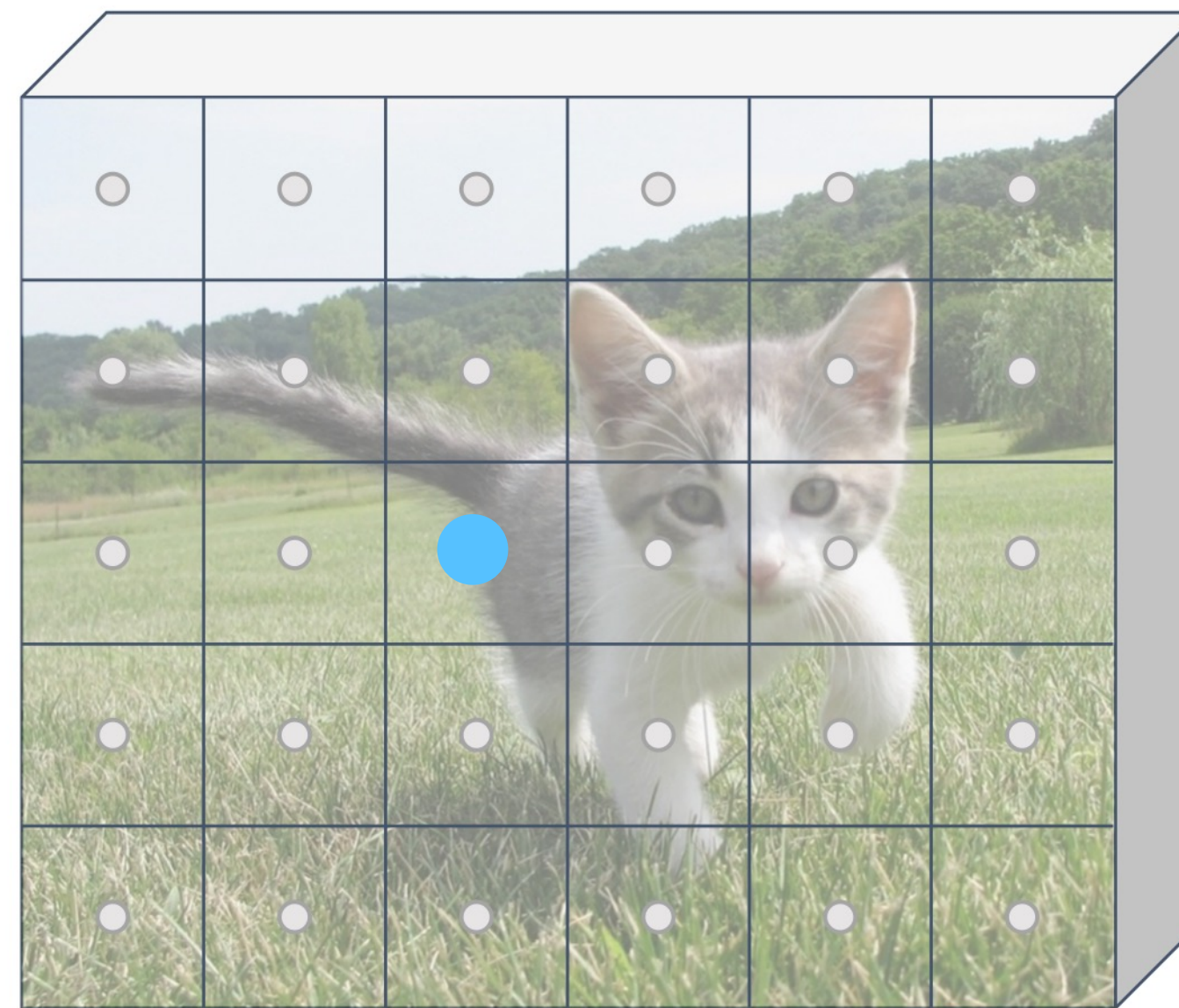


Image features
(e.g. 512 x 5 x 6)

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



Anchor is object?
 $2K \times 5 \times 6$

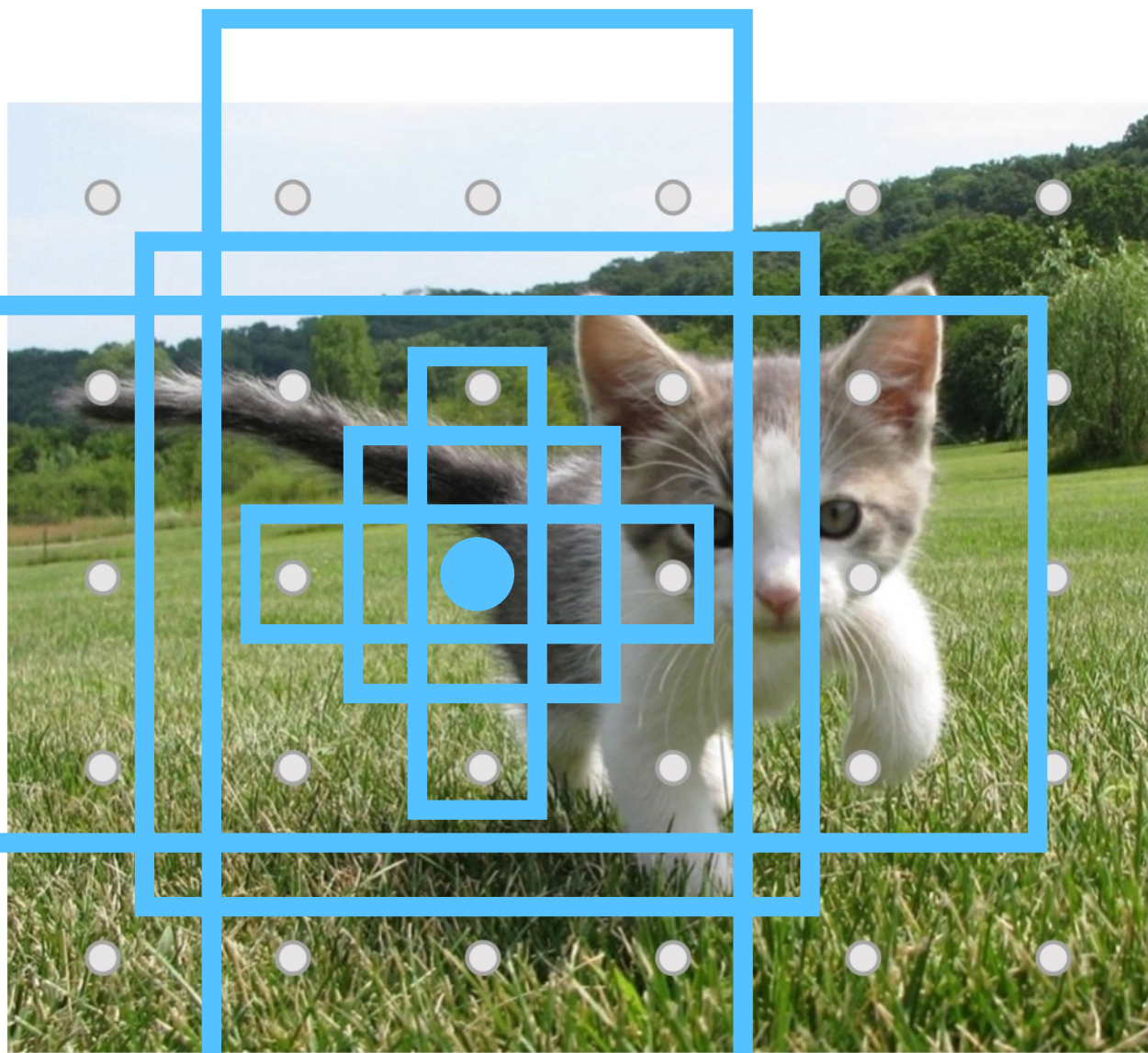


Anchor transforms
 $4K \times 5 \times 6$



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

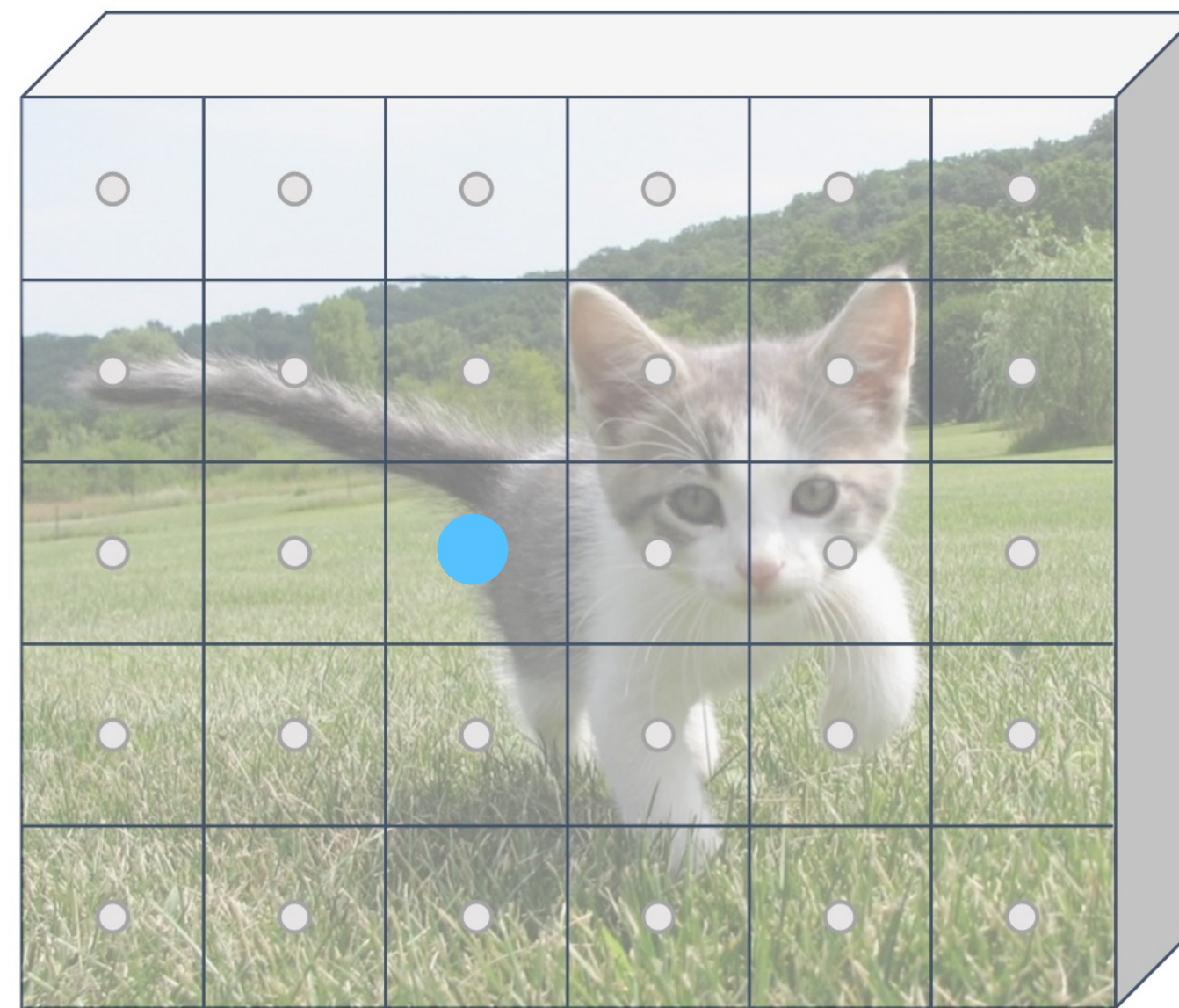


Image features
(e.g. 512 x 5 x 6)

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)



Anchor is object?
2K x 5 x 6



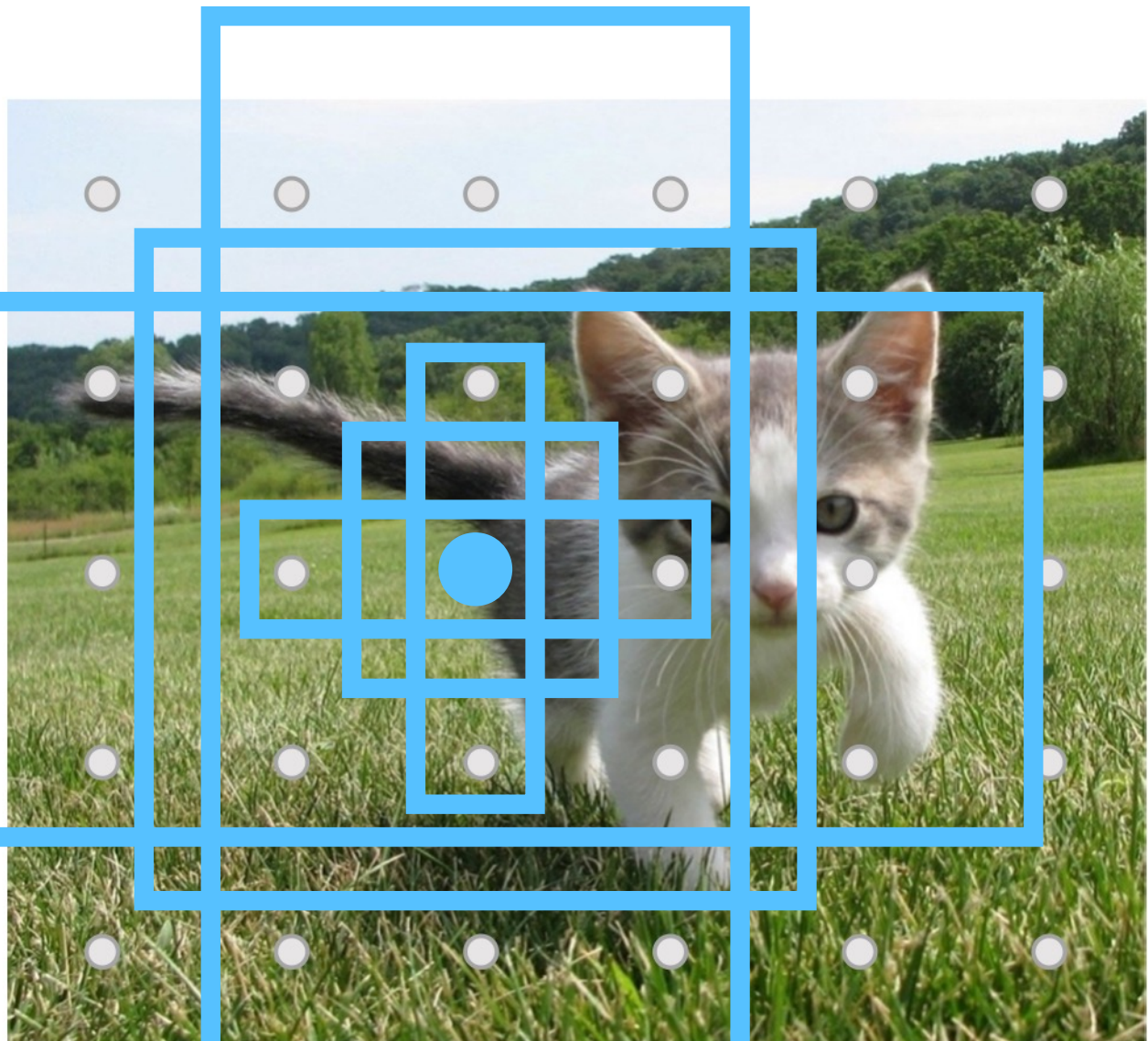
Anchor transforms
4K x 5 x 6

During training, supervised positive / negative anchors and box transforms like R-CNN



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

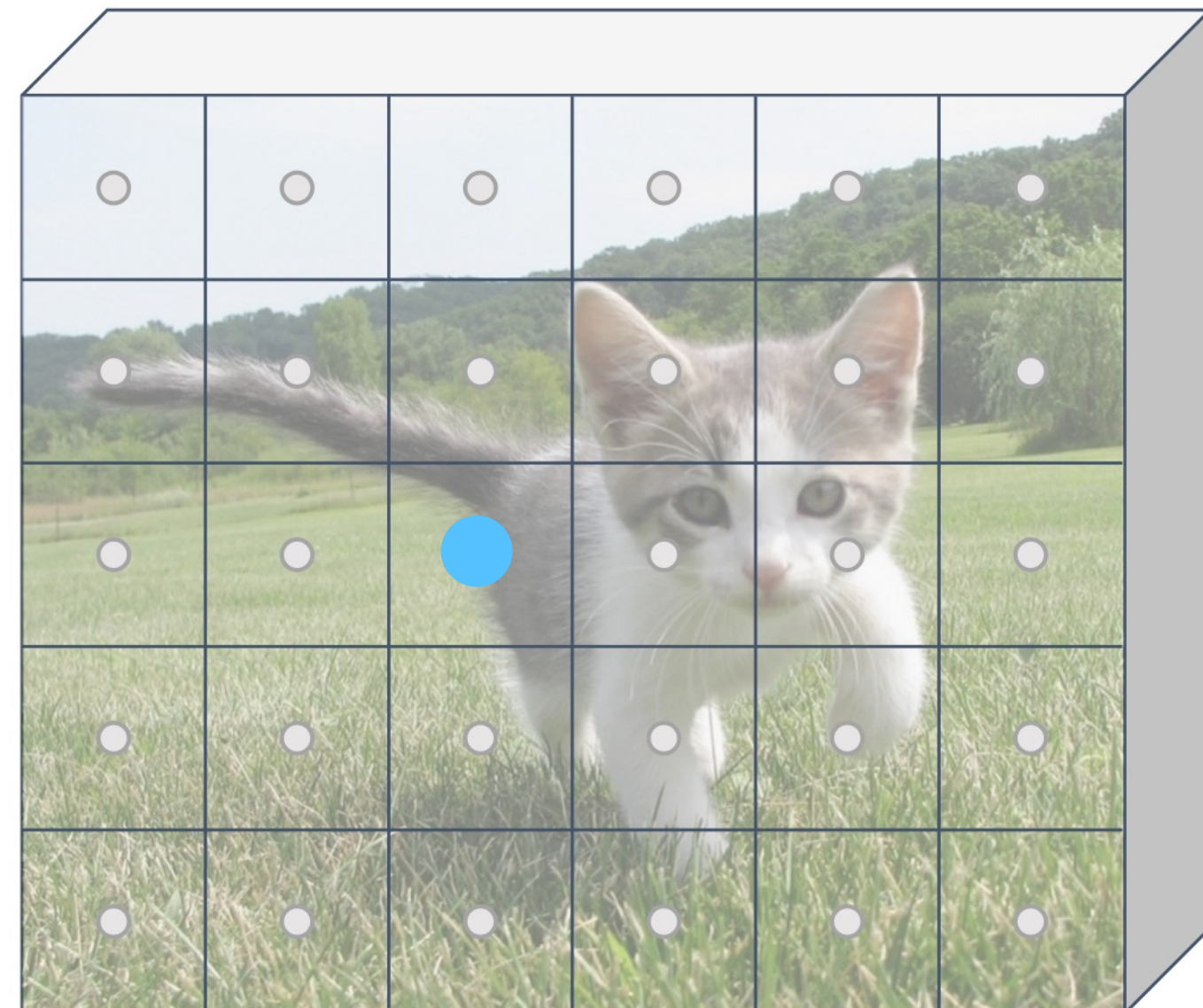
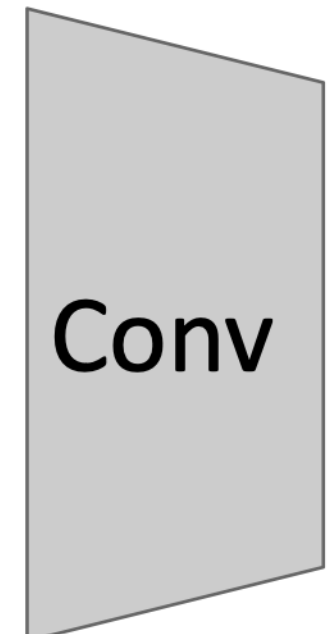


Image features
(e.g. 512 x 5 x 6)

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)



Anchor is object?
2K x 5 x 6



Anchor transforms
4K x 5 x 6

Positive anchors: ≥ 0.7 IoU with some GT box (plus highest IoU to each GT)

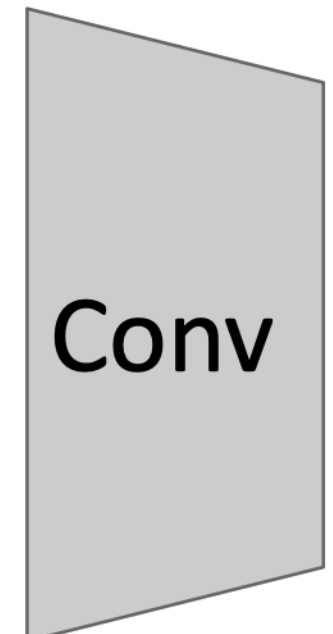
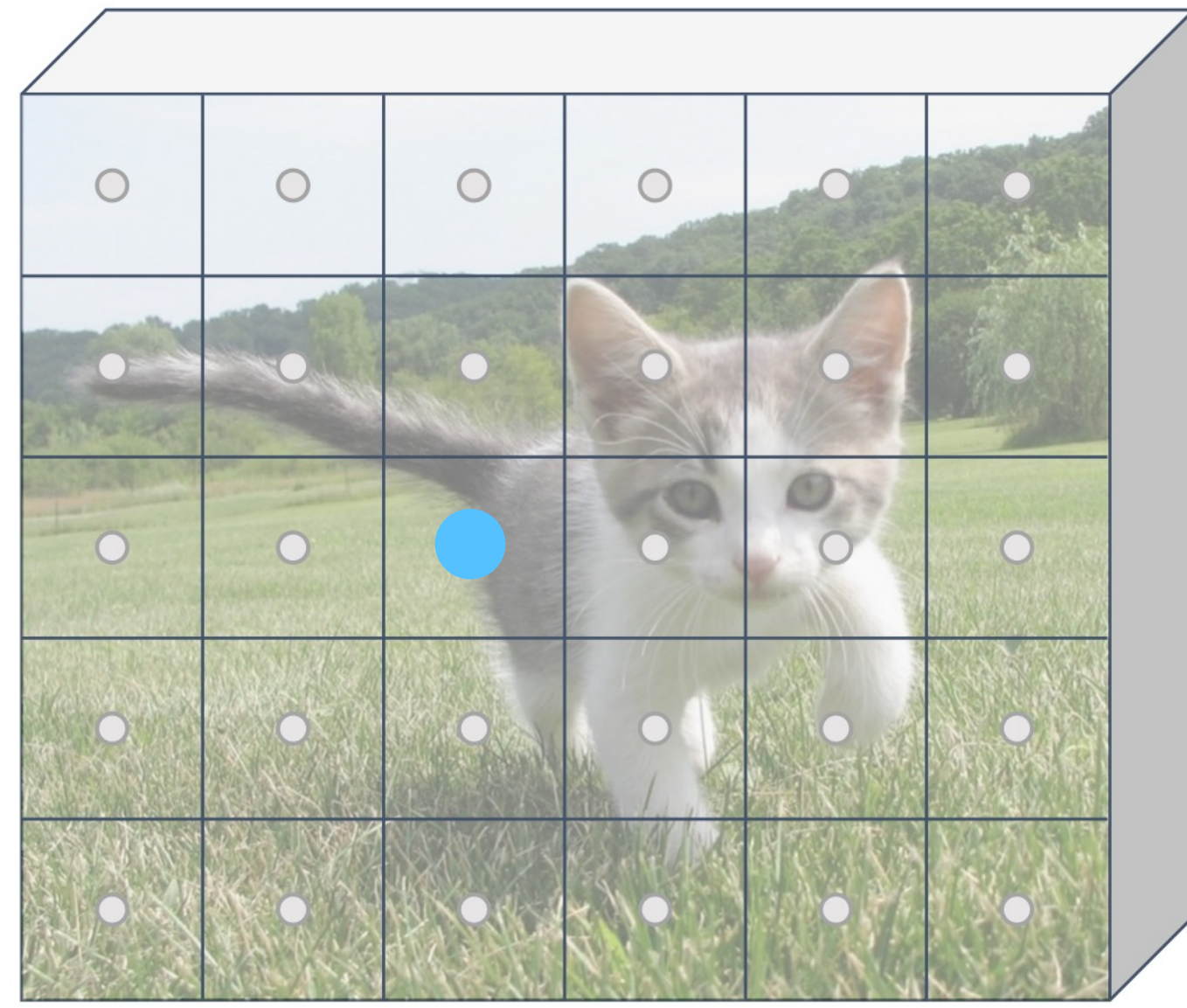
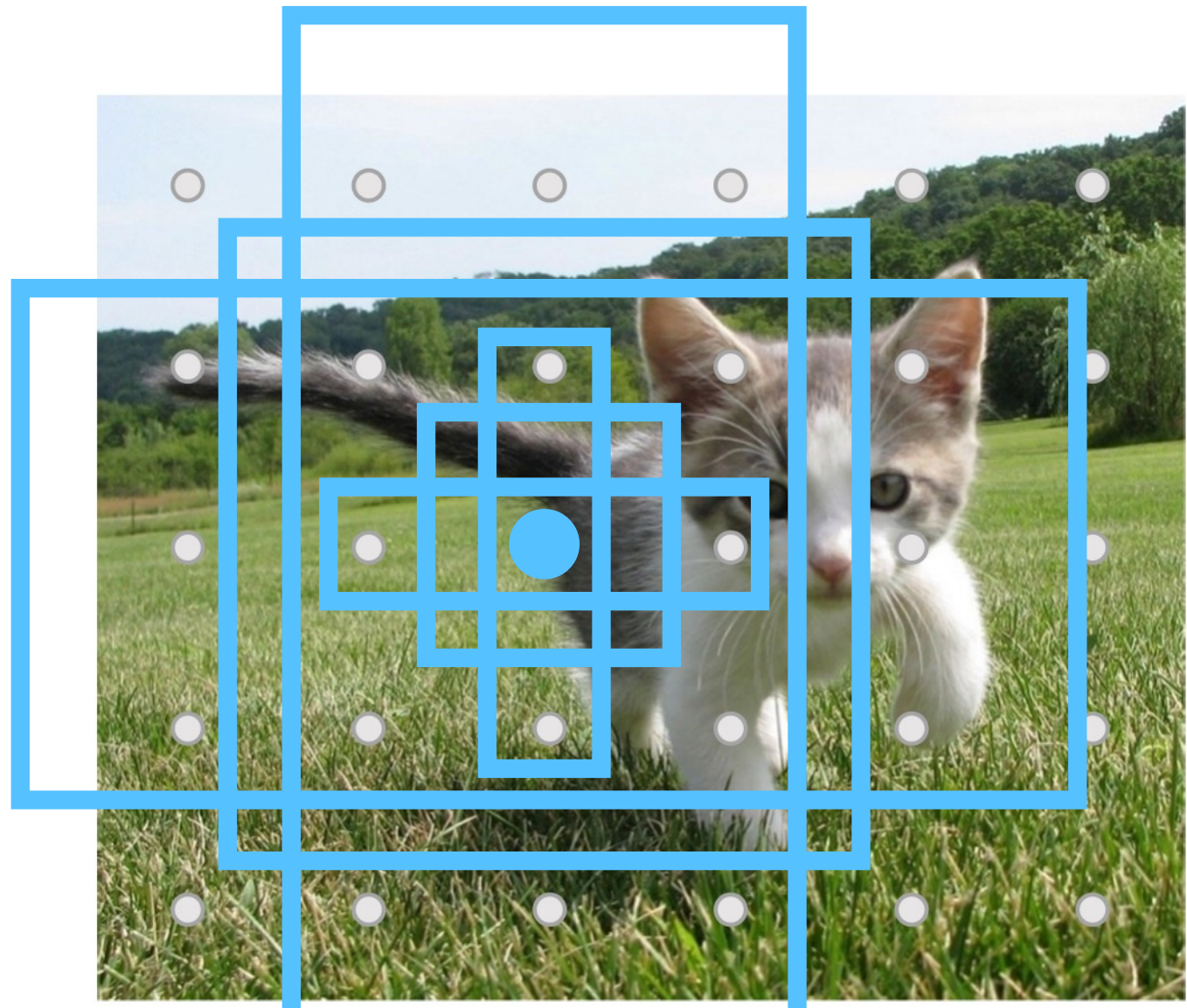


Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image

Each feature corresponds to a point in the input

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)



Anchor is object?
2K x 5 x 6

Anchor transforms
4K x 5 x 6

Input Image
(e.g. 3 x 640 x 480)

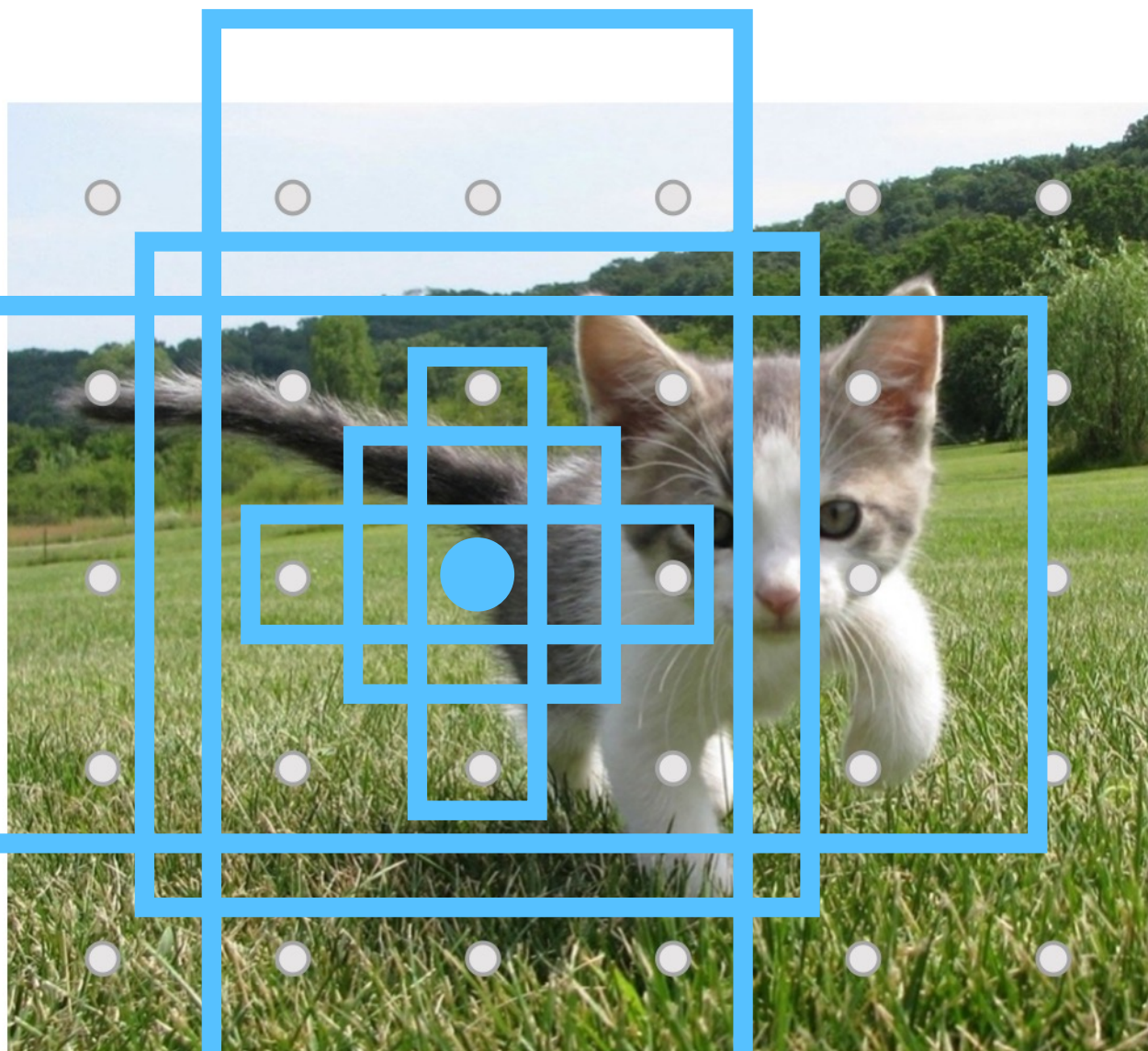
Image features
(e.g. 512 x 5 x 6)

Negative anchors: < 0.3 IoU with all GT boxes. Don't supervised transforms for negative boxes.



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

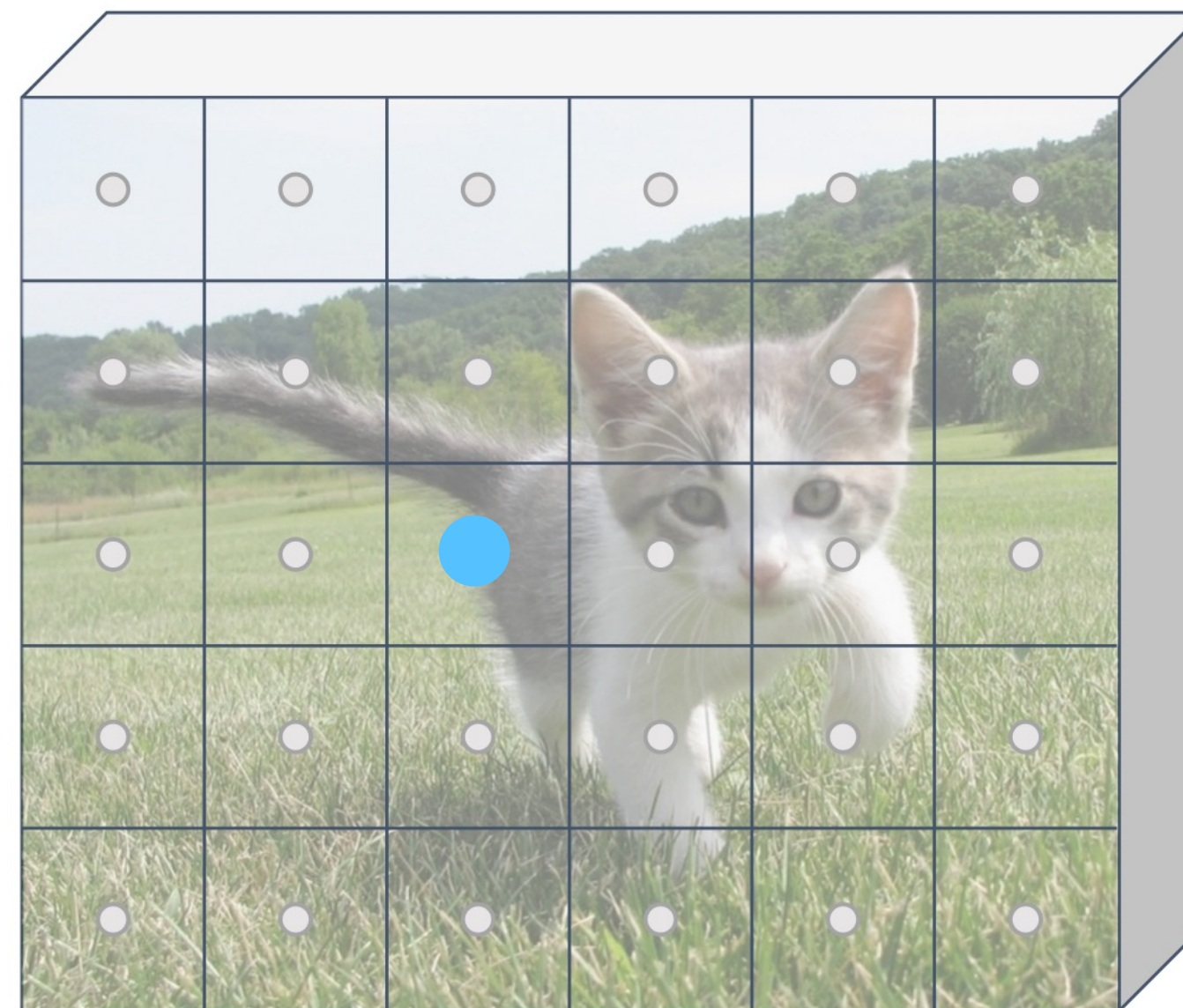


Image features
(e.g. 512 x 5 x 6)

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here K = 6)



Anchor is object?
2K x 5 x 6



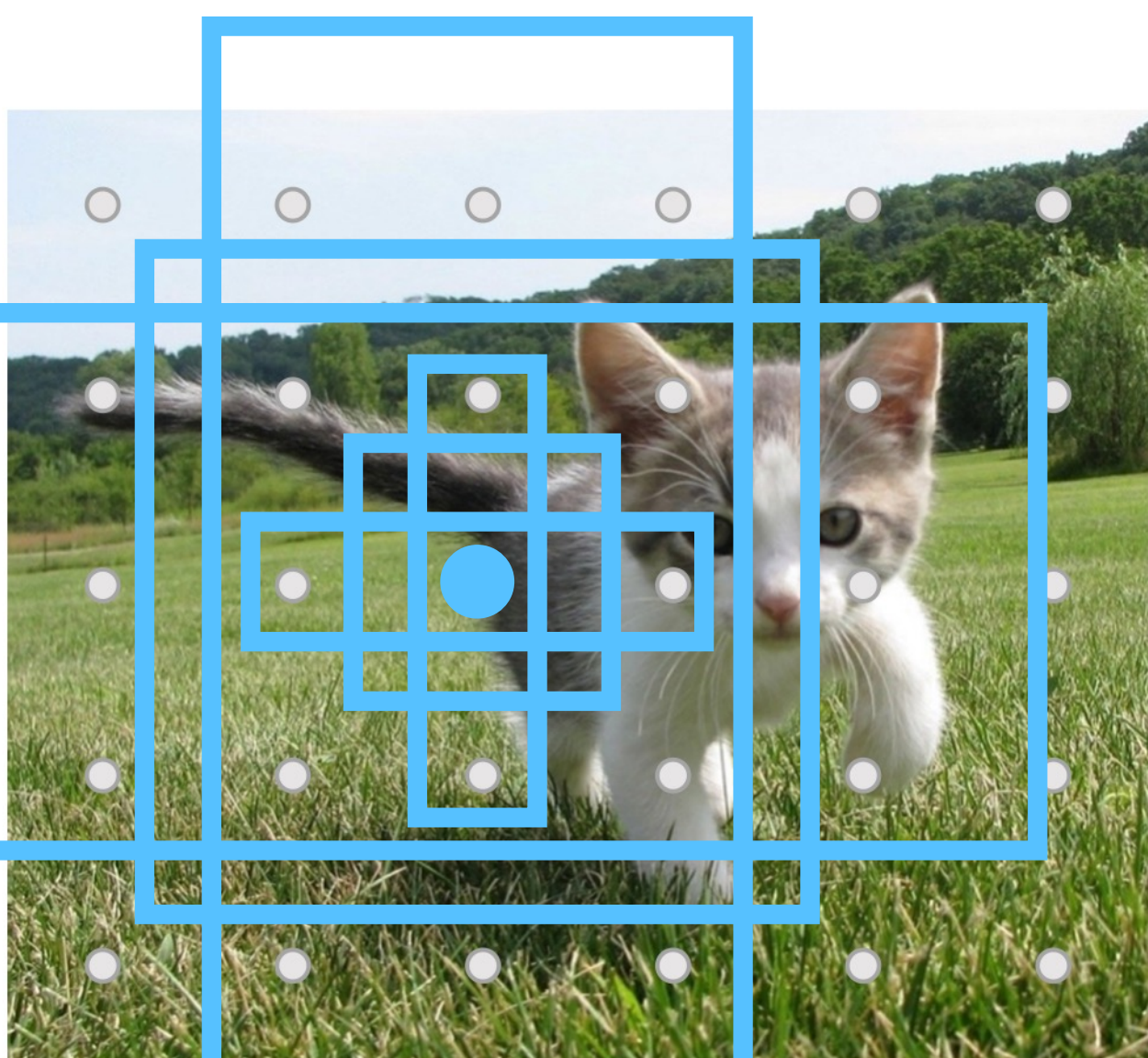
Anchor transforms
4K x 5 x 6

Neutral anchors: between 0.3 and 0.7 IoU with all GT boxes; ignored during training



Region Proposal Network (RPN)

Run backbone CNN to get features aligned to input image



Input Image
(e.g. 3 x 640 x 480)



Each feature corresponds to a point in the input

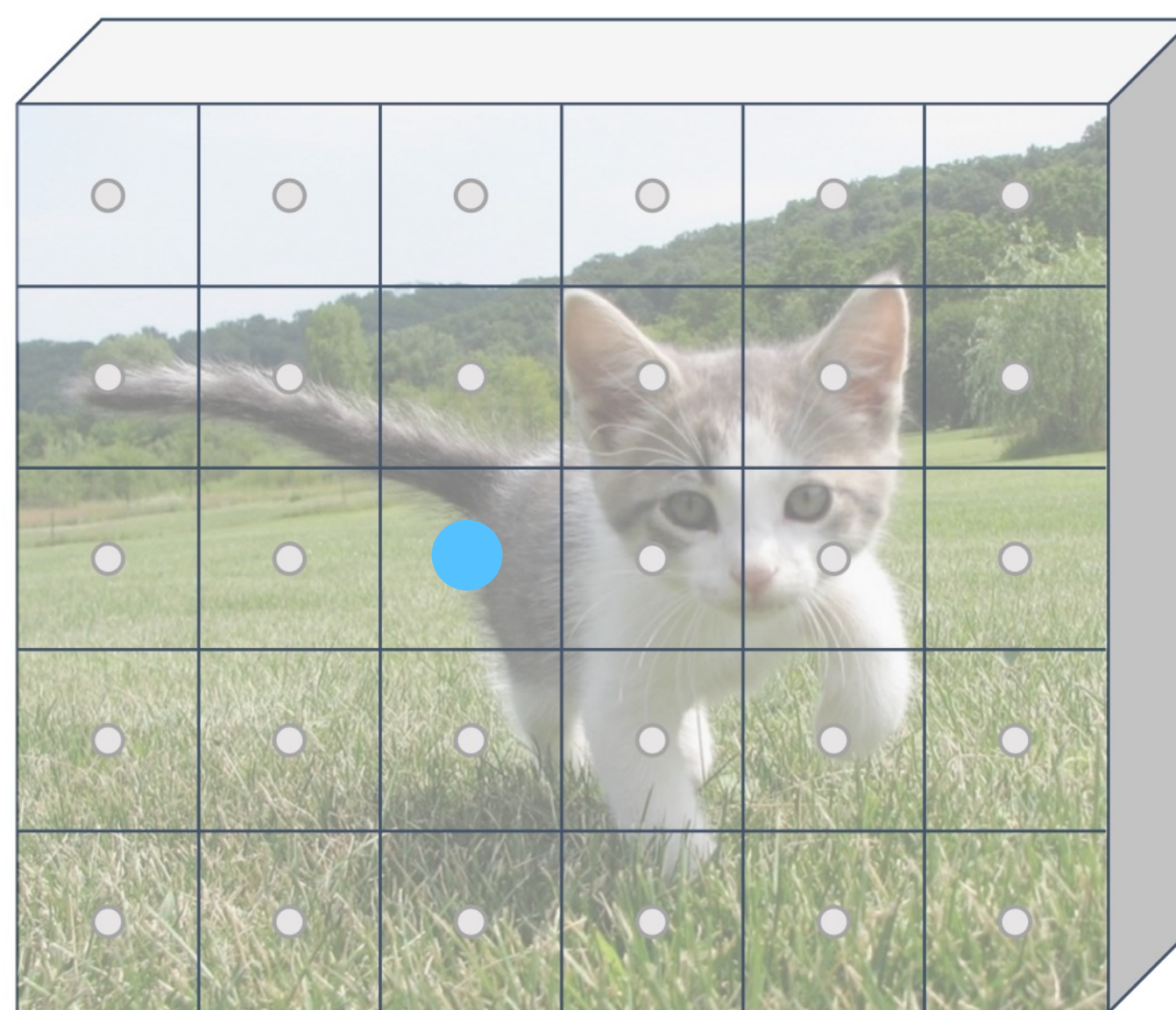


Image features
(e.g. 512 x 5 x 6)

In practice: Rather than using one anchor per point, instead consider K different anchors with different size and scale (here $K = 6$)



Anchor is object?
 $2K \times 5 \times 6$



Anchor transforms
 $4K \times 5 \times 6$

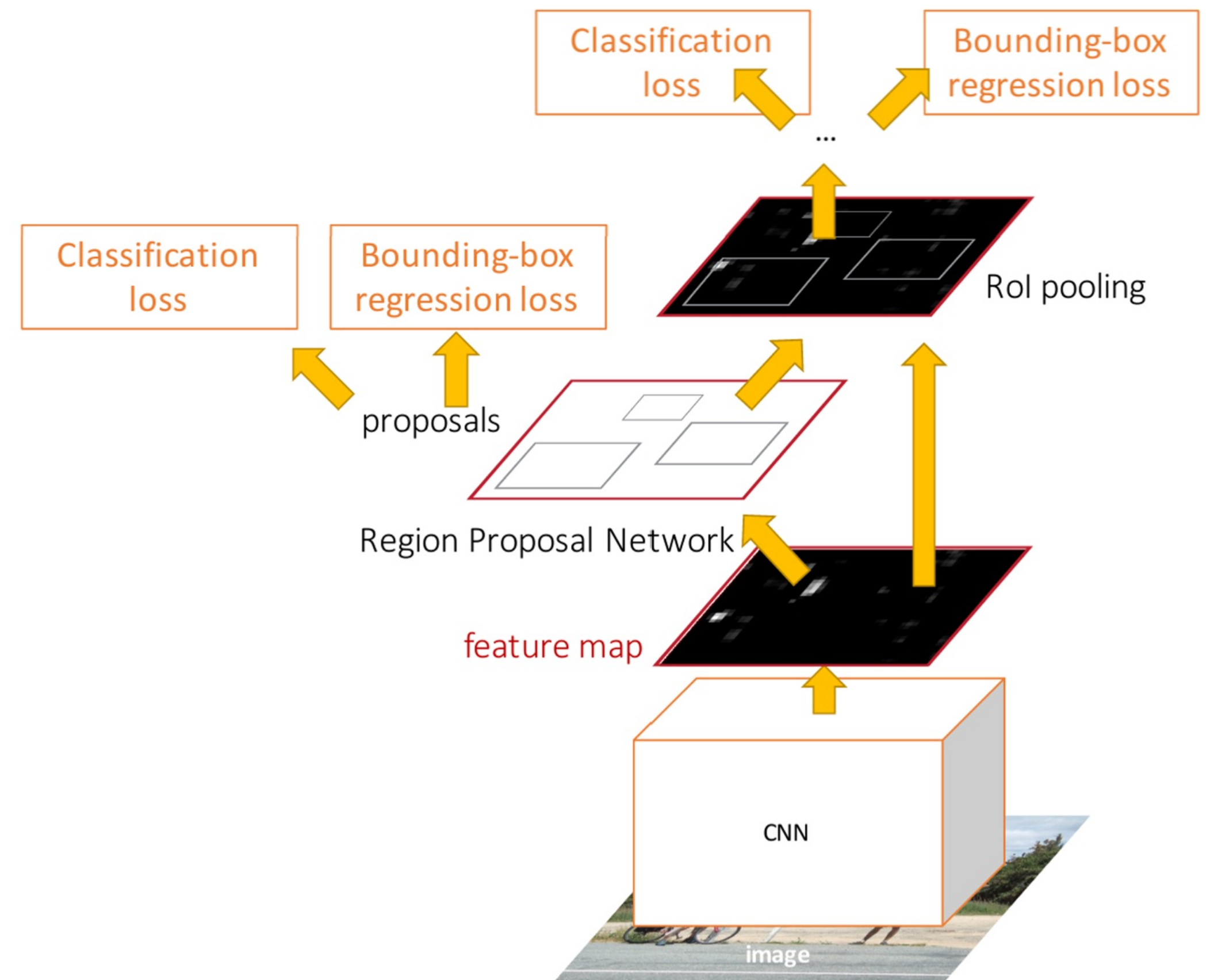
At test-time, sort all $K \cdot 5 \cdot 6$ boxes by their positive score, take top 300 as our region proposals



Faster R-CNN: Learnable Region Proposals

Jointly train four losses:

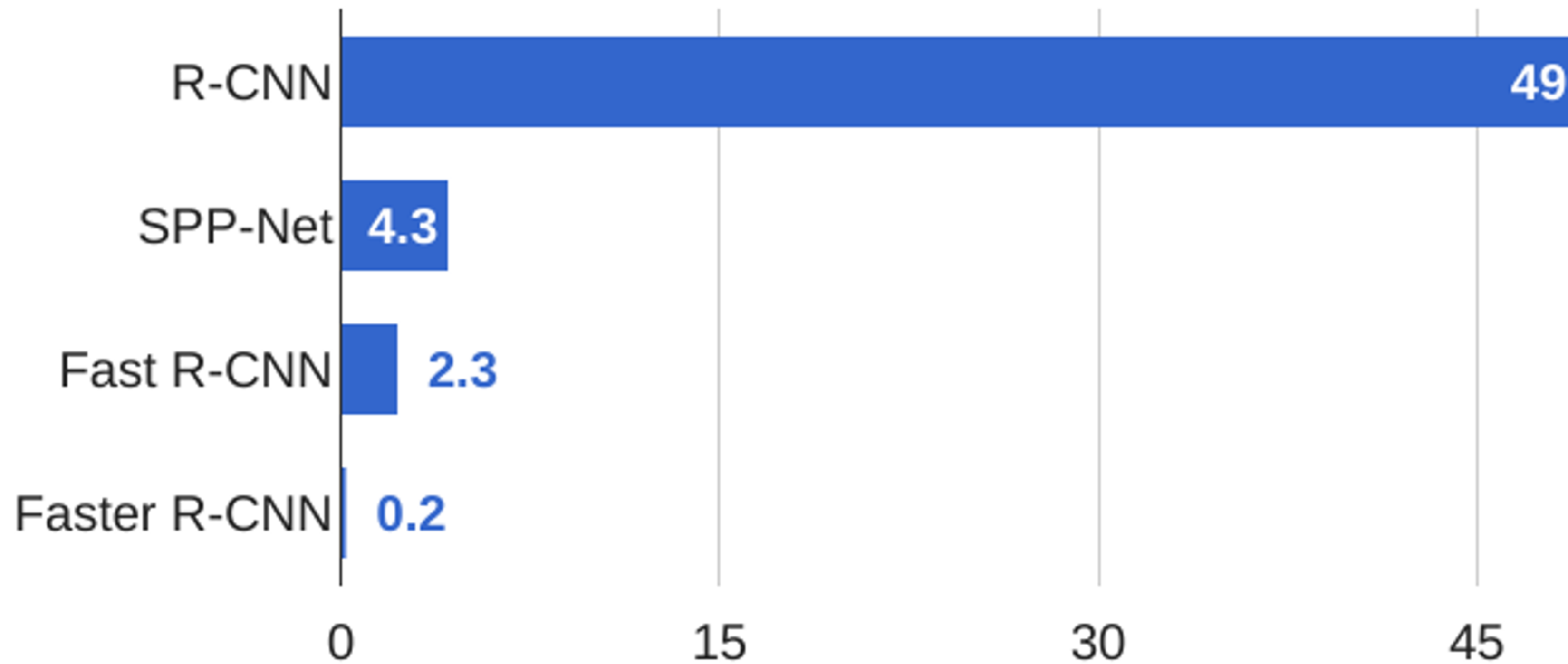
1. **RPN classification:** anchor box is object / not an object
2. **RPN regression:** predict transform from anchor box to proposal box
3. **Object classification:** classify proposals as background / object class
4. **Object regression:** predict transform from proposal box to object box





Faster R-CNN: Learnable Region Proposals

R-CNN Test-Time Speed (s)





Extend Faster R-CNN to Image Segmentation: Mask R-CNN

Classification



“Chocolate Pretzels”

← No spatial extent →

Semantic

Segmentation



Chocolate Pretzels,

Shelf

← No objects, just pixels →

Object

Detection



Flipz, Hershey's, Keese's

← Multiple objects →

Instance

Segmentation





Extend Faster R-CNN to Instance Segmentation: Mask R-CNN

Instance Segmentation

Detect all objects in the image and identify the pixels that belong to each object (Only things!)

Approach

Perform object detection then predict a segmentation mask for each object detected!

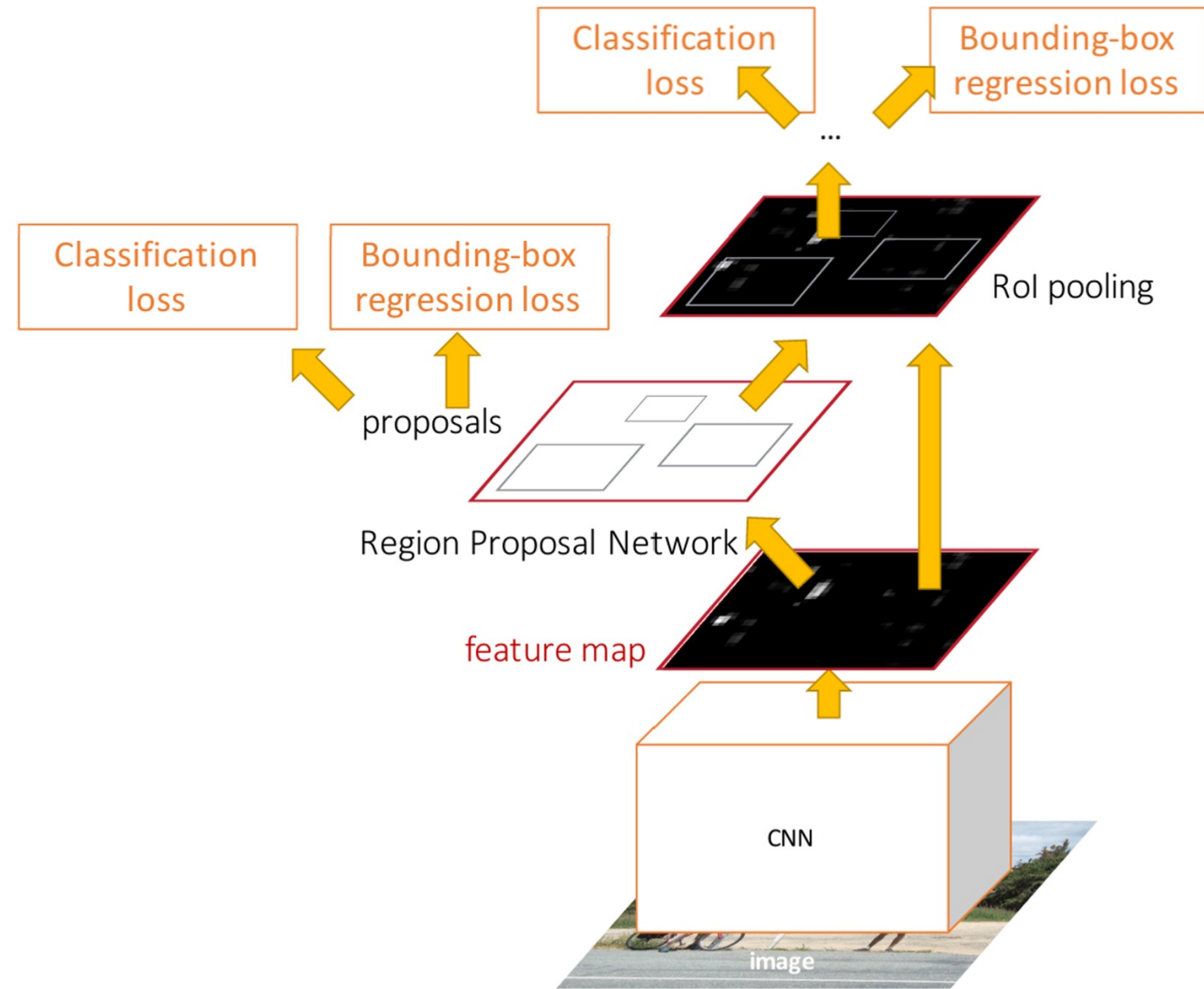




Extend Faster R-CNN into Mask R-CNN

Faster R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 1. **Object classification:** classify proposals
 2. **Object regression:** predict transform from proposal box to object box



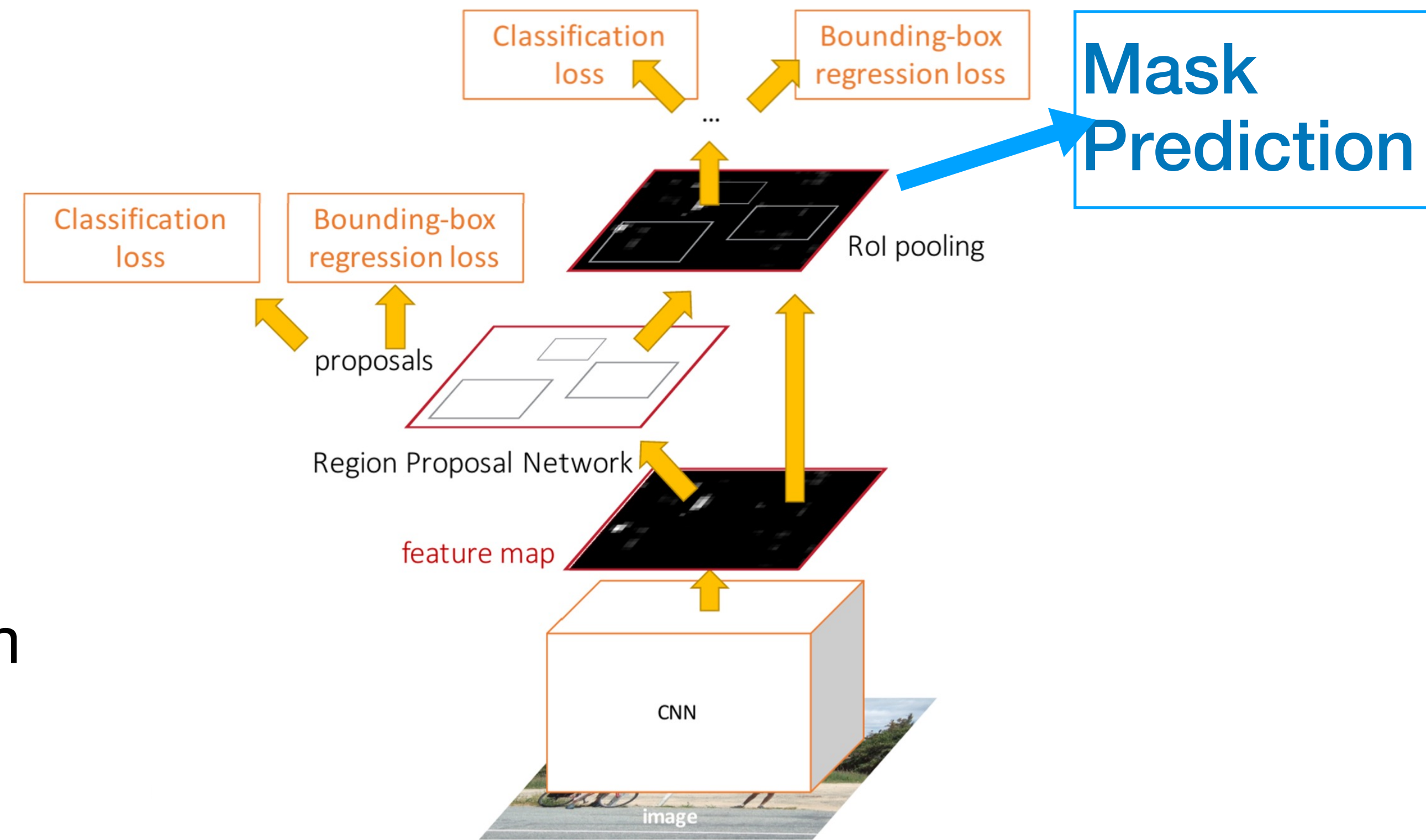


Extend Faster R-CNN into Mask R-CNN

Faster R-CNN

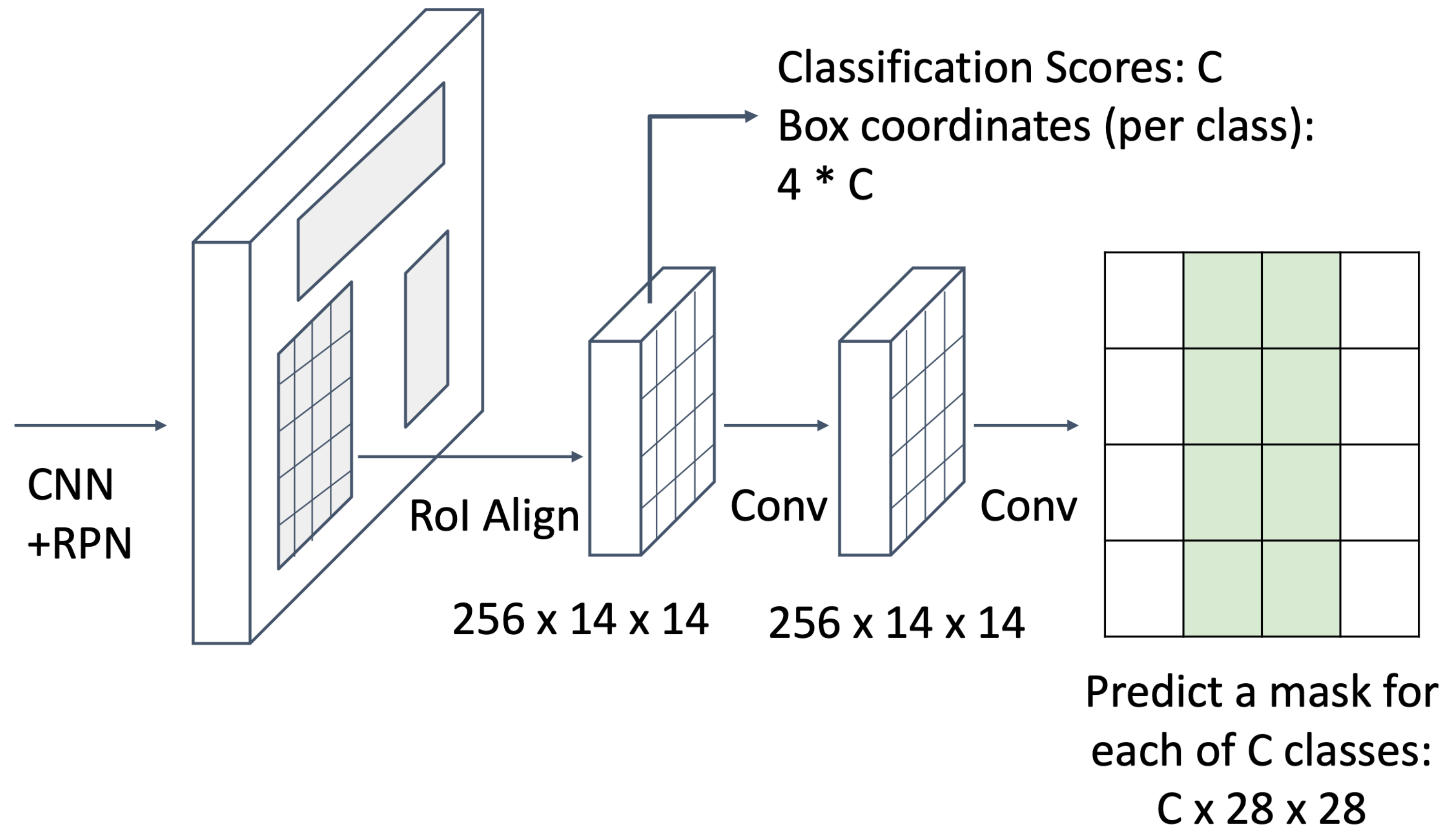
Mask R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region





Mask R-CNN





Mask R-CNN: Very Good Results!

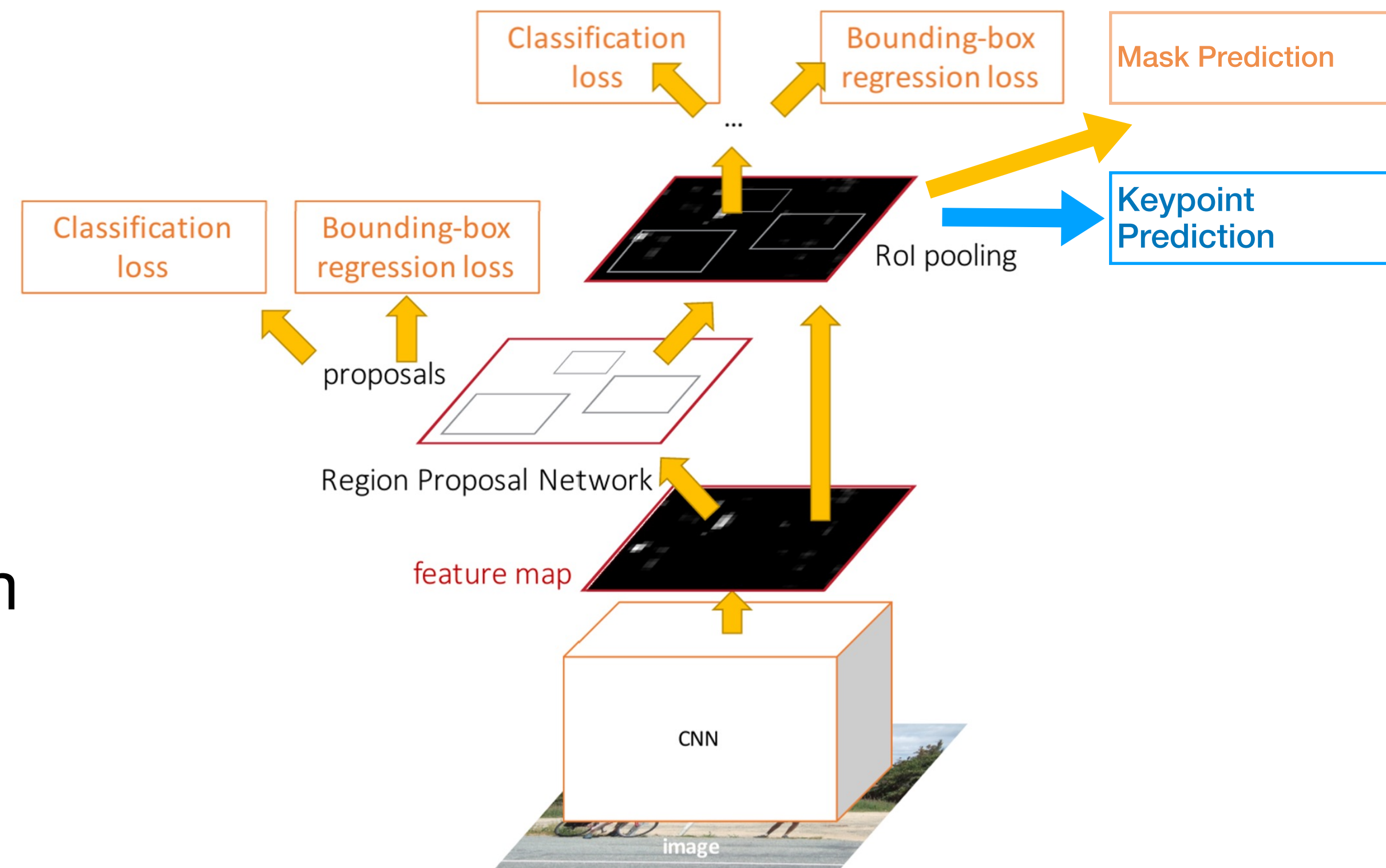




Mask R-CNN for Human Pose Estimation

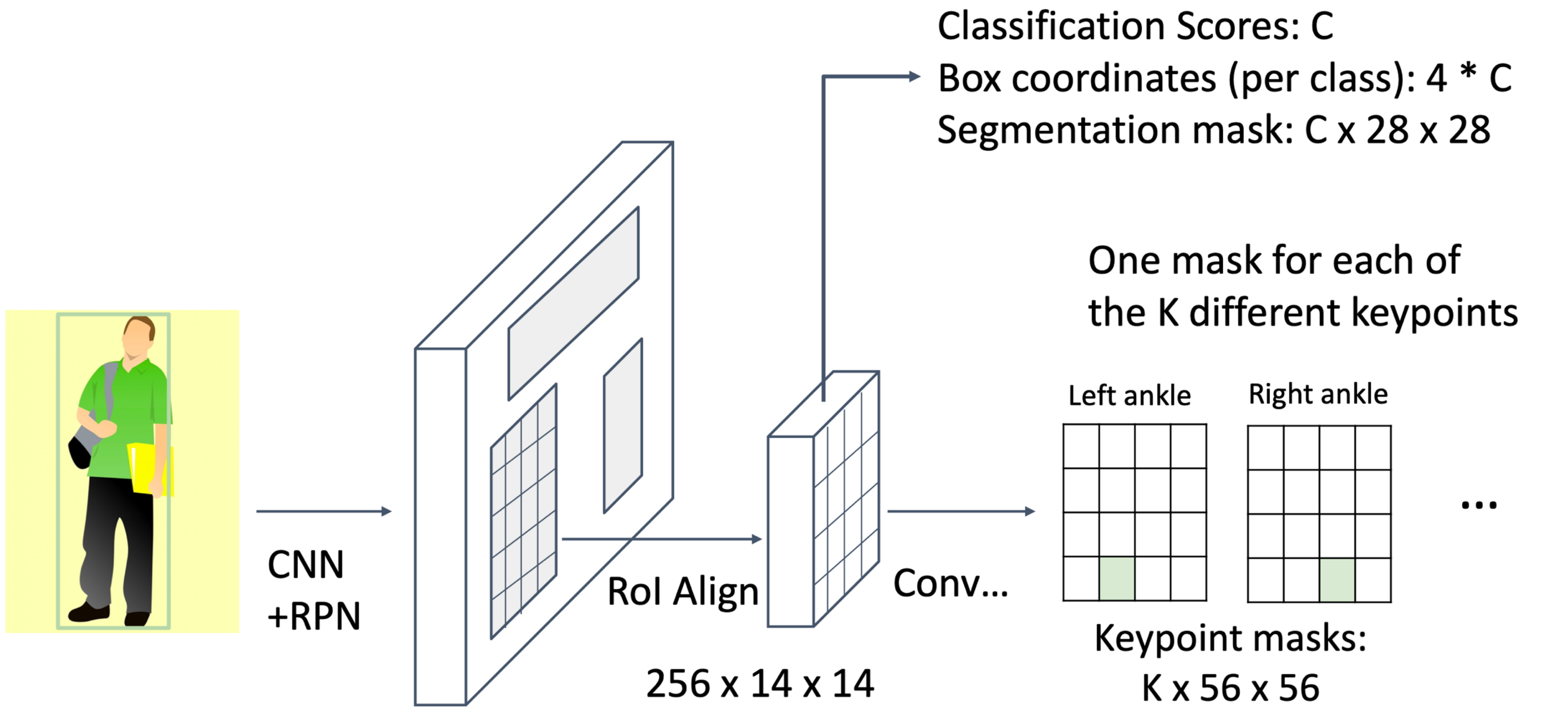
Mask R-CNN

1. **Feature Extraction** at the image-level
2. **Regions of Interest** proposal from feature map
3. **In Parallel**
 - a. **Object Classification:** classify proposals
 - b. **Object Regression:** predict transform from proposal box to object box
 - c. **Mask Prediction:** predict a binary mask for every region
 - d. **Keypoint Prediction:** predict binary mask for human key points





Mask R-CNN for Human Pose Estimation



Ground-truth has one "pixel" turned on per keypoint. Train with softmax loss



Mask R-CNN for Human Pose Estimation





Two Stage vs One Stage Detectors

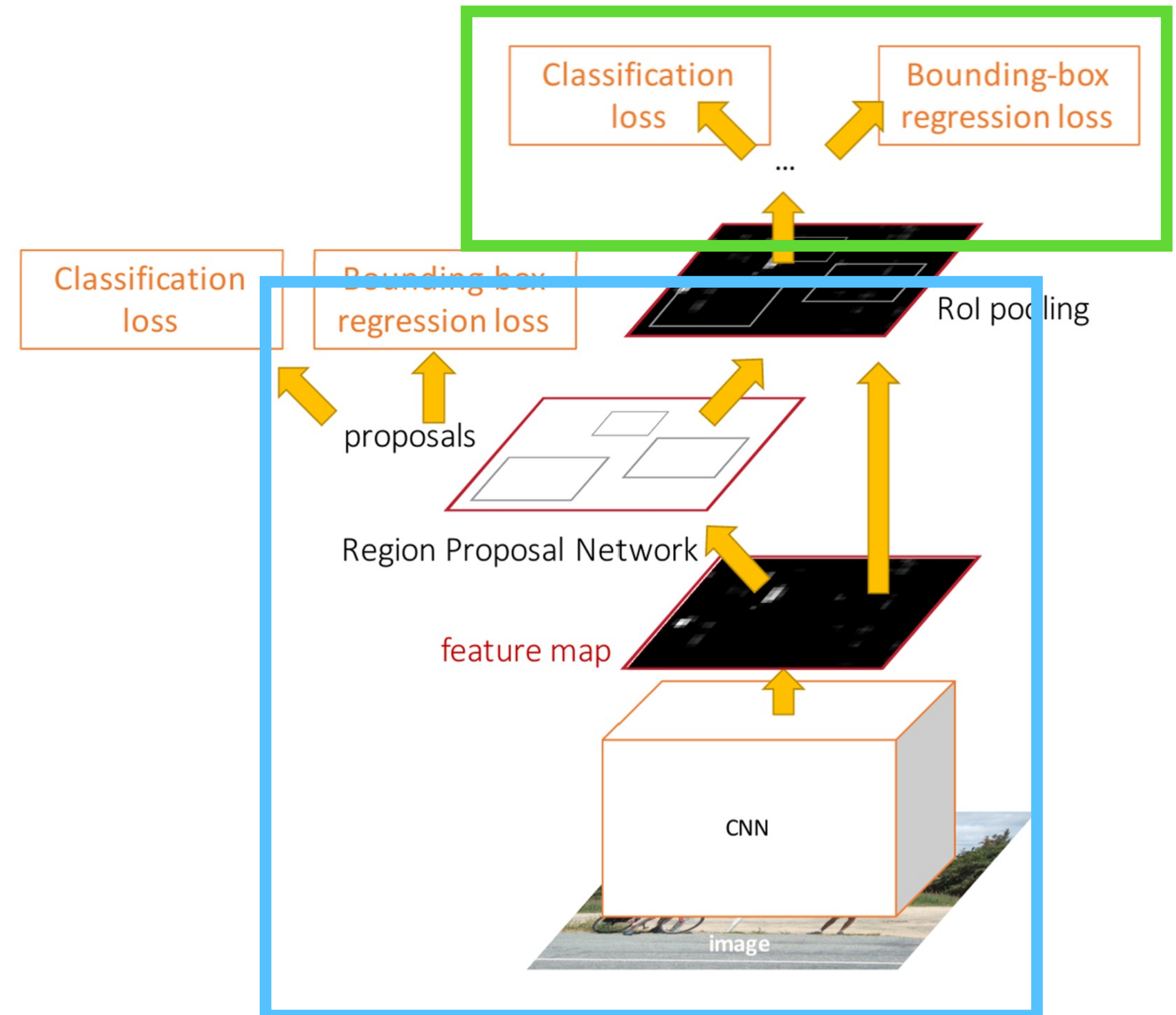
Faster R-CNN is a two-stage object detector

First stage. Run once per image

- Backbone Network

Second stage. Run once per region

- Crop features: RoI pool / align
- Predict Object Class
- Prediction bbox offset



Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks"



DEEP ROB

Lecture 9
Object Detection
University of Michigan | Department of Robotics